

Global SMEFT fits guided by ML

1st COMETA General Meeting

29/02/24

Izmir

JHEP 03 (2023) 033

2211.02058

Jaco ter Hoeve

VU Amsterdam & Nikhef Theory Group



Outline

Mapping the SMEFT with SMEFiT

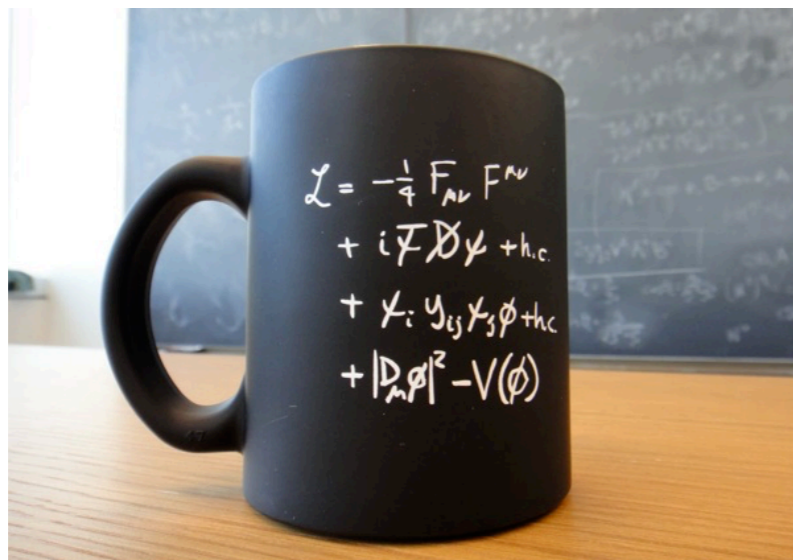
ML assisted observables: the ML4EFT framework

Case study: ML observables in the Higgs and top sector

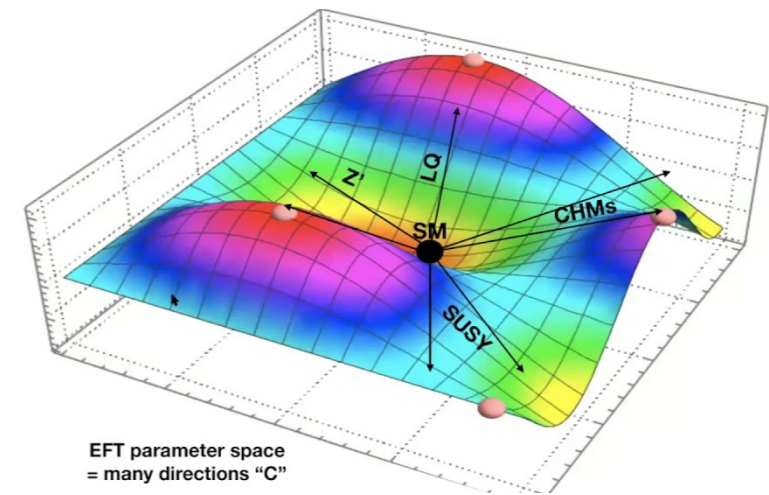
A combined ML4EFT + SMEFiT analysis

The Standard Model as an EFT

$$\mathcal{L}_{\text{SMEFT}} = \mathcal{L}_{\text{SM}} + \sum_i^{N_{d5}} \frac{c_i}{\Lambda} \mathcal{O}_i^{(5)} + \sum_i^{N_{d6}} \frac{c_i}{\Lambda^2} \mathcal{O}_i^{(6)} + \sum_i^{N_{d7}} \frac{c_i}{\Lambda^3} \mathcal{O}_i^{(7)} + \sum_i^{N_{d8}} \frac{b_i}{\Lambda^4} \mathcal{O}_i^{(8)} + \dots$$



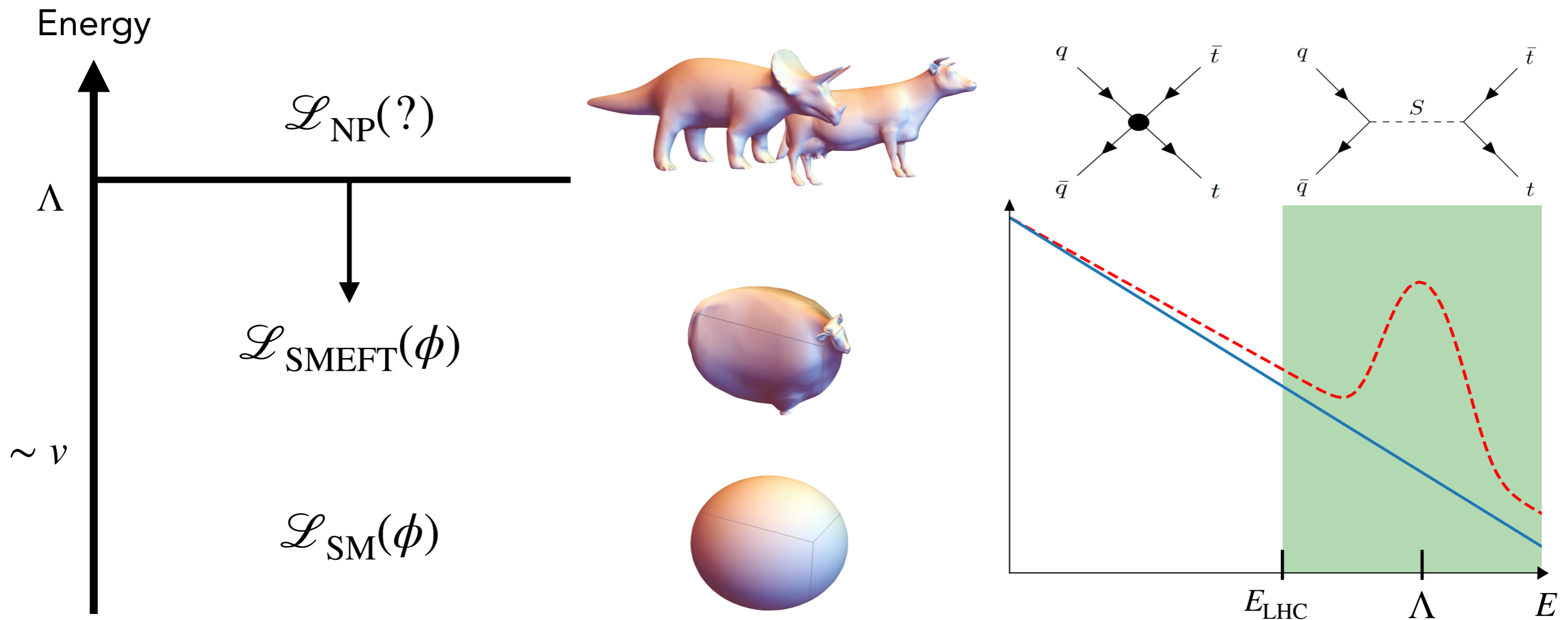
+



- ▶ **Low energy limit** of generic UV-complete theories at high energies
- ▶ Assumes the **SM field content and symmetries**
- ▶ **Complete basis** at any given mass dimension

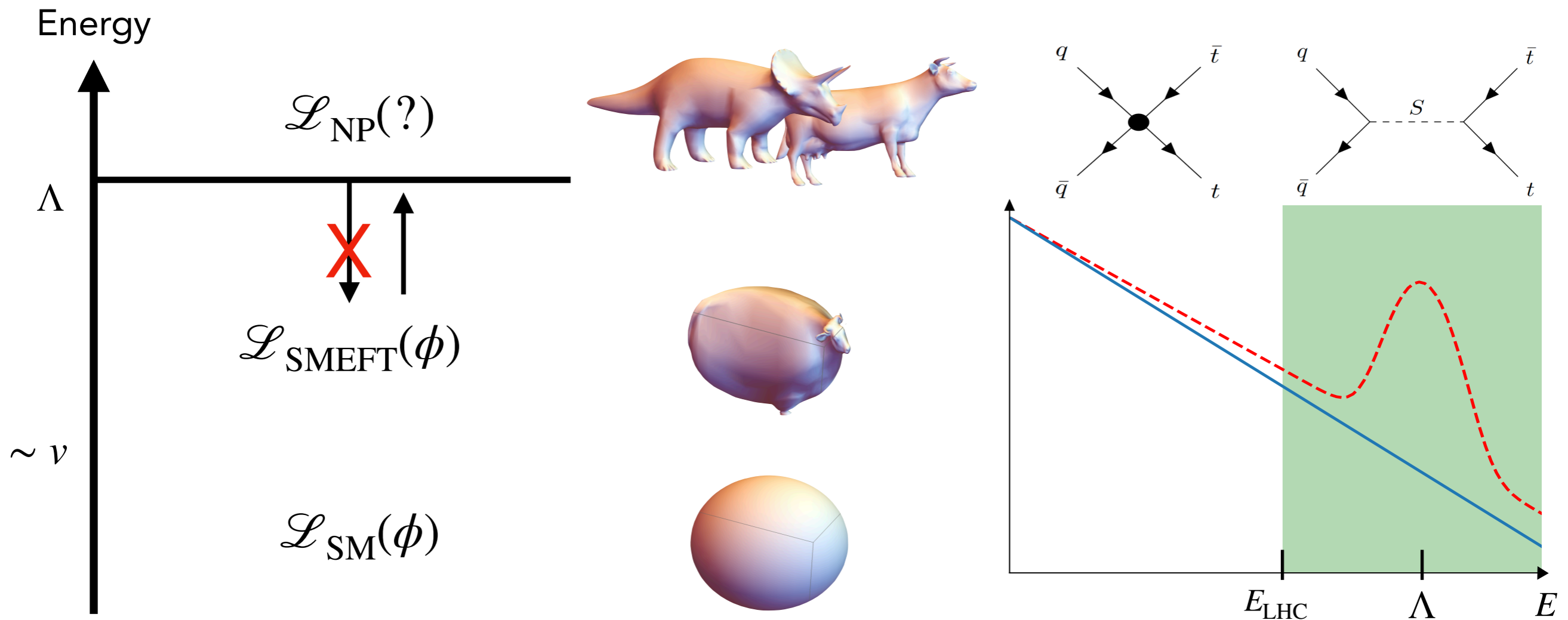
The Standard Model as an EFT

$$\mathcal{L}_{\text{SMEFT}} = \mathcal{L}_{\text{SM}} + \sum_i^{N_{d5}} \frac{c_i}{\Lambda} \mathcal{O}_i^{(5)} + \sum_i^{N_{d6}} \frac{c_i}{\Lambda^2} \mathcal{O}_i^{(6)} + \sum_i^{N_{d7}} \frac{c_i}{\Lambda^3} \mathcal{O}_i^{(7)} + \sum_i^{N_{d8}} \frac{b_i}{\Lambda^4} \mathcal{O}_i^{(8)} + \dots$$



The Standard Model as an EFT

$$\mathcal{L}_{\text{SMEFT}} = \mathcal{L}_{\text{SM}} + \sum_i^{N_{d5}} \frac{c_i}{\Lambda} \mathcal{O}_i^{(5)} + \sum_i^{N_{d6}} \frac{c_i}{\Lambda^2} \mathcal{O}_i^{(6)} + \sum_i^{N_{d7}} \frac{c_i}{\Lambda^3} \mathcal{O}_i^{(7)} + \sum_i^{N_{d8}} \frac{b_i}{\Lambda^4} \mathcal{O}_i^{(8)} + \dots$$



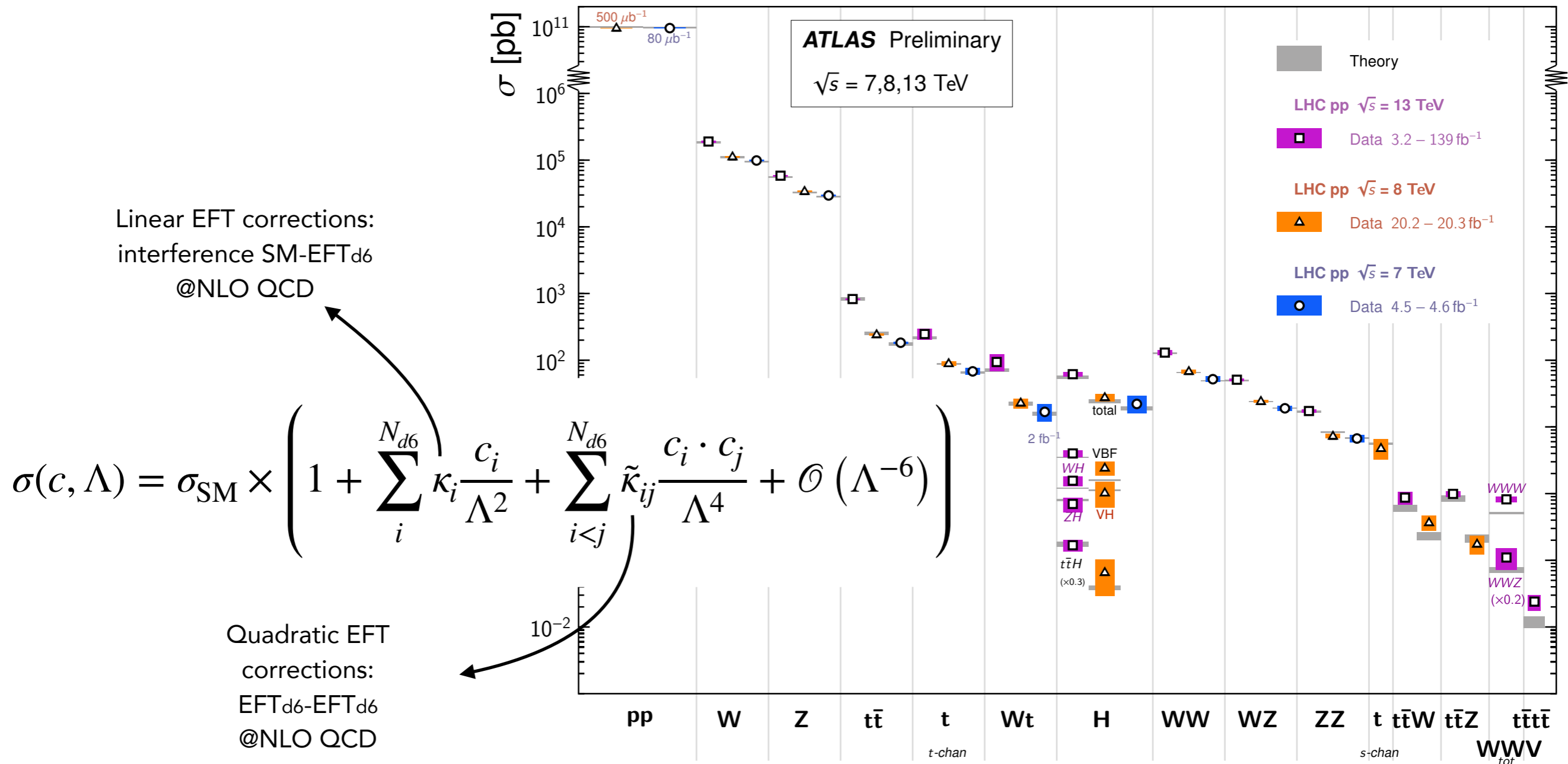
LHC observables in the SMEFT

Idea: parameterise (differential) cross-sections in terms of higher dimensional operators

[ATL-PHYS-PUB-2022-009]

Standard Model Total Production Cross Section Measurements

Status: February 2022



LHC observables in the SMEFT

From (differential) cross sections ...

$$\sigma_{\text{SMEFT}}(c, \Lambda) = \sigma_{\text{SM}} \times \left(1 + \sum_i^{N_{d6}} \kappa_i \frac{c_i}{\Lambda^2} + \sum_{i < j}^{N_{d6}} \tilde{\kappa}_{ij} \frac{c_i \cdot c_j}{\Lambda^4} + \mathcal{O}(\Lambda^{-6}) \right)$$

Linear EFT corrections:
interference SM-EFT_{d6}
@NLO QCD

Quadratic EFT
corrections:
EFT_{d6}-EFT_{d6}
@NLO QCD

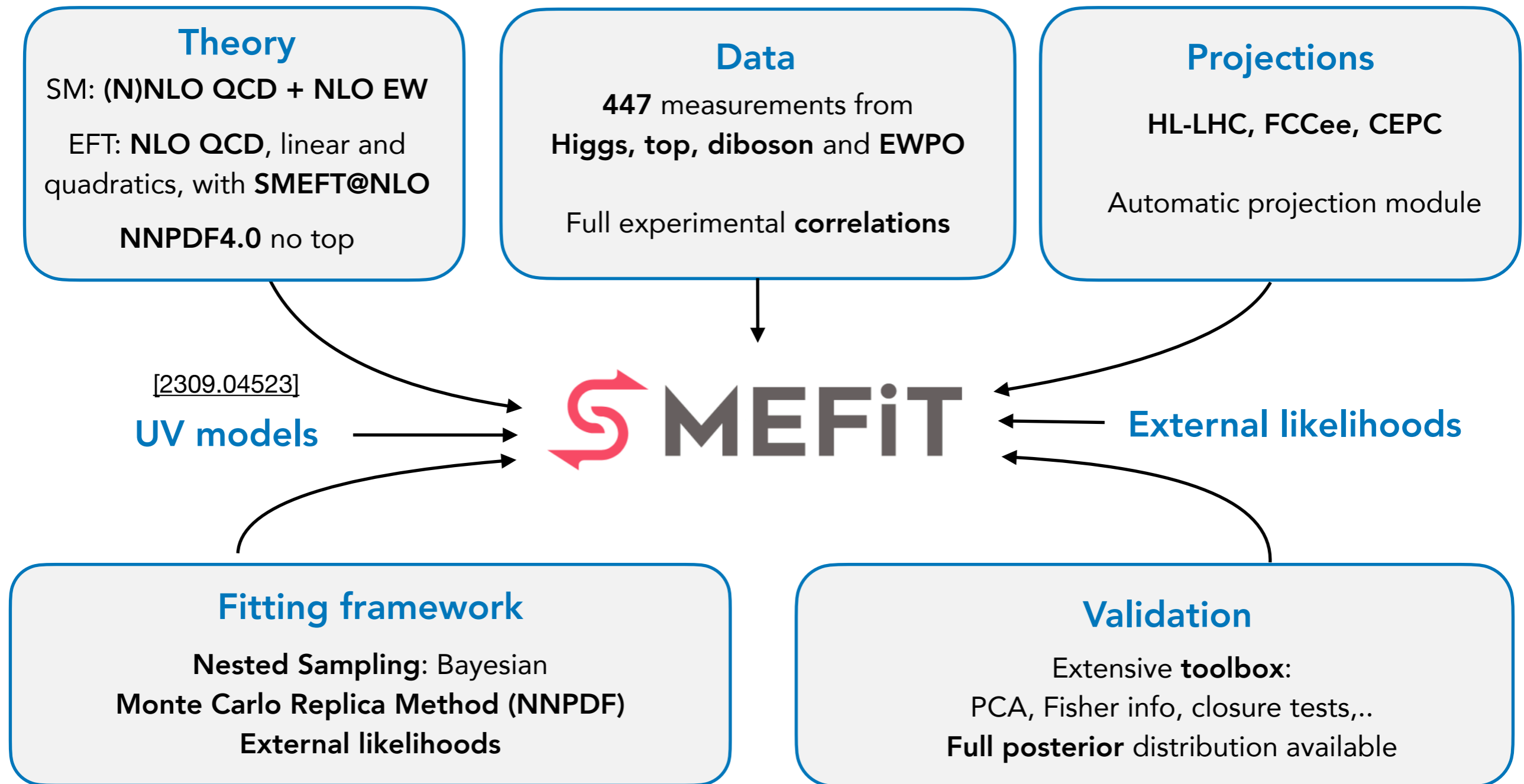
To a combined likelihood ready for optimisation ...

$$-2 \log \mathcal{L} = \frac{1}{n_{\text{dat}}} \sum_{i,j=1}^{n_{\text{dat}}} \left(\sigma_{i,\text{SMEFT}}(c) - \sigma_{i,\text{exp}} \right) (\text{cov}^{-1})_{ij} \left(\sigma_{j,\text{SMEFT}}(c) - \sigma_{j,\text{exp}} \right)$$

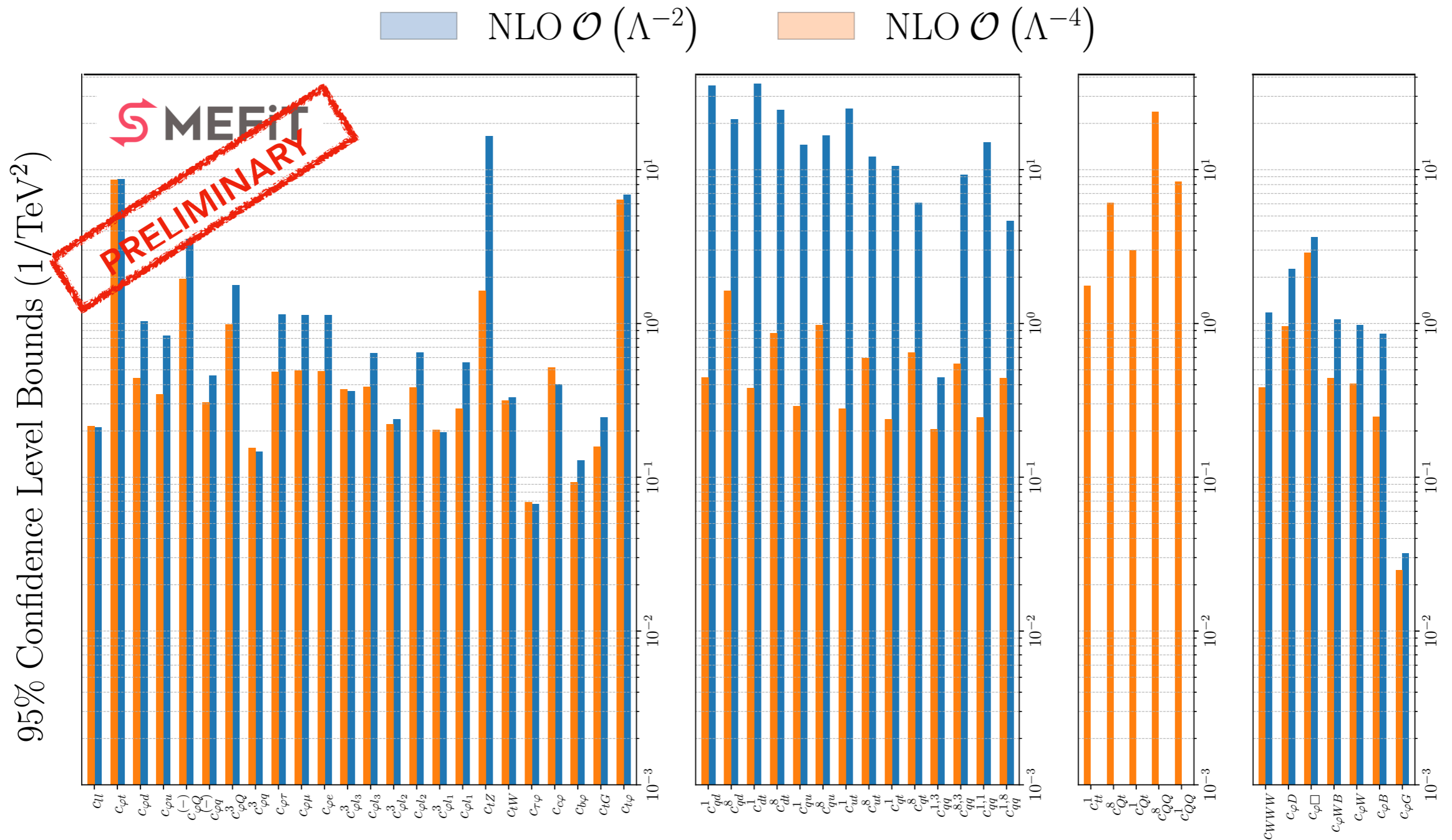
Theory (pdf + scale) and experimental uncertainties (stat + systematics): $\text{cov}^{(\text{tot})}_{ij} = \text{cov}^{(\text{th})}_{ij} + \text{cov}^{(\text{exp})}_{ij}$

The SMEFiT framework

[2302.06660]



Result: state of the art



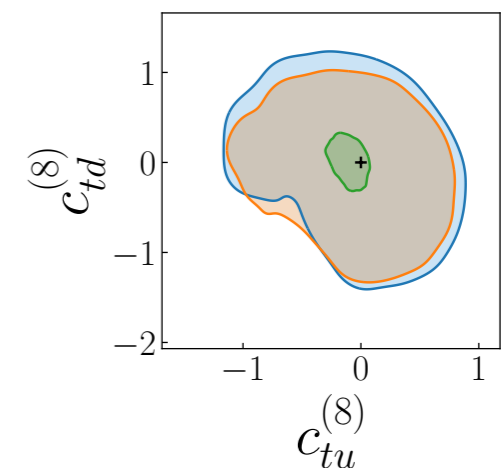
Global EFT fits include data on **top quark, Higgs, and gauge boson, EWPO**, both inclusive and differential measurements for a total of 447 measurements

But can we do (even) better?

- State of the art global efforts **reinterpret** "SM measurements" in an EFT context
- Which measurement is the most **sensitive** to EFT parameters?
 - Inclusive, single to multi-differential (which variables)
 - Binned or unbinned, which binning?

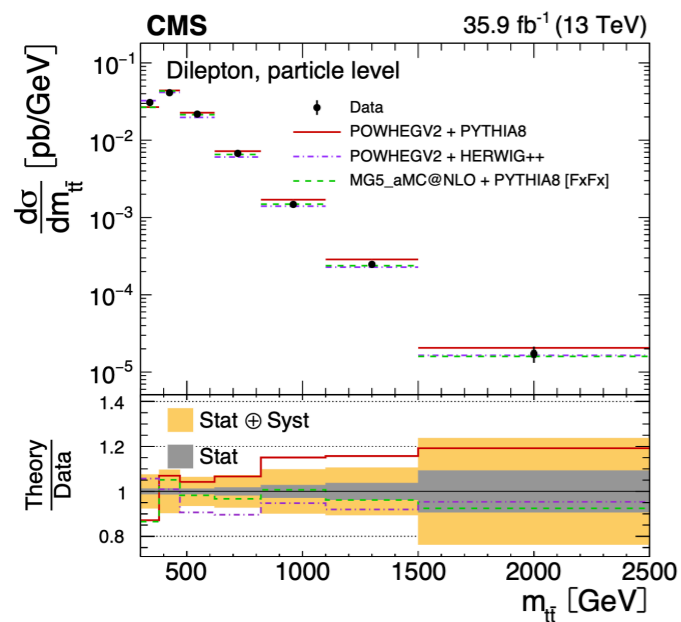
Framework needed to integrate **unbinned multivariate observables** into global SMEFT fits

- **Optimal** bounds on the EFT parameters
- Useful **diagnosis tool** to assess information loss

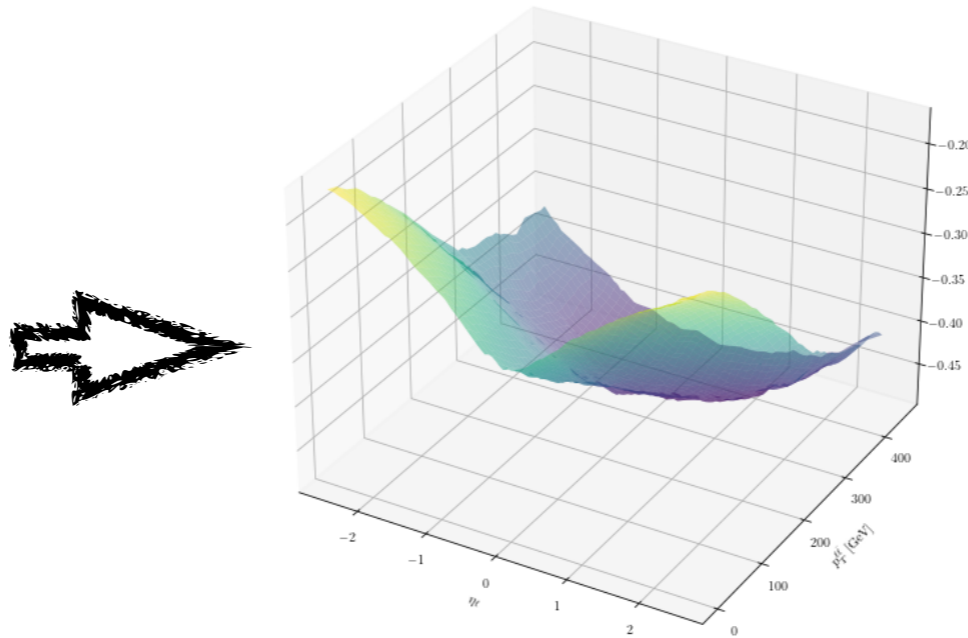


But can we do (even) better?

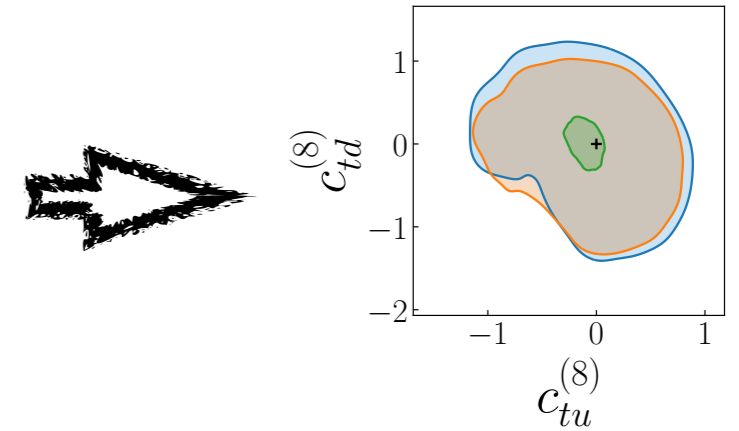
- State of the art global efforts reinterpret "SM measurements" in an EFT context



Binned, univariate

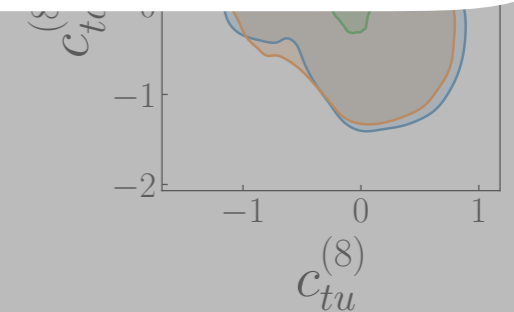


Unbinned, multivariate



Optimal inference

- Optimal bounds on the EFT parameters
- Useful diagnosis tool to assess information loss



In the remainder of this talk



Mapping the SMEFT with SMEFiT

Unbinned multivariate observables: the ML4EFT framework

Case study: ML observables in the Higgs and top sector

A combined ML4EFT + SMEFiT analysis

ML4EFT

[2211.02058] R. Gomez Ambrosio, JtH, M. Madigan, J. Rojo, V.Sanz

<https://lhcfiteknikhef.github.io/ML4EFT>

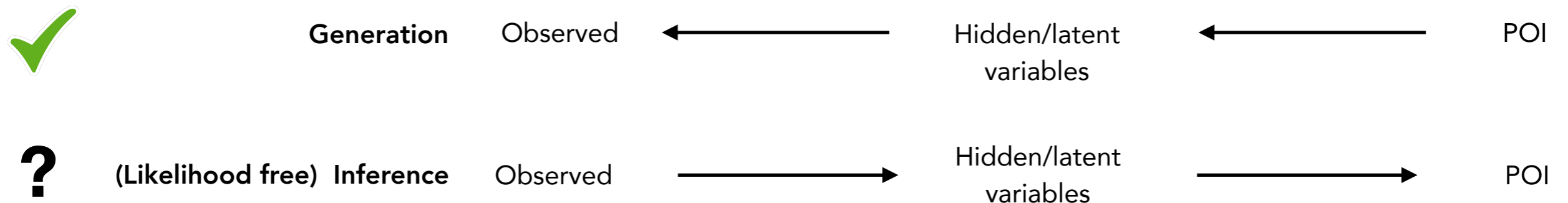
Where theory and experiment meet

We are progressively moving through the simulation chain (latent space)

$$p(x|c) \sim \int dz_{\text{det}} dz_{\text{shower}} dz_{\text{parton}} p(x|z_{\text{det}}) p(z_{\text{det}}|z_{\text{shower}}) p(z_{\text{shower}}|z_{\text{parton}}) p(z_{\text{parton}}|c)$$



Unbinned unfolding, Omnifold [1911.09107]



Where theory and experiment meet

We are progressively moving through the simulation chain (latent space)

$$p(x|c) \sim \int dz_{\text{det}} dz_{\text{shower}} dz_{\text{parton}} p(x|z_{\text{det}}) p(z_{\text{det}}|z_{\text{shower}}) p(z_{\text{shower}}|z_{\text{parton}}) p(z_{\text{parton}}|c)$$



Unbinned unfolding, Omnifold [1911.09107]



Generation Observed ← Hidden/latent variables ← POI

? (Likelihood free) Inference Observed → Hidden/latent variables → POI

Likelihood free inference

- Starting from two balanced datasets \mathcal{D}_{SM} and \mathcal{D}_{EFT} drawn from $f(\mathbf{x} | \text{SM})$ and $f(\mathbf{x} | \text{EFT})$, we minimise e.g. the cross-entropy loss

$$L[g(\mathbf{x})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{EFT}}} w_e \log(1 - g(\mathbf{x}_e)) - \frac{1}{N} \sum_{\mathcal{D}_{\text{SM}}} w_e \log g(\mathbf{x}_e)$$

Event weights
{ $m_{\bar{t}\bar{t}}, \eta_l, \Delta\phi, \dots$ }

- The learned decision boundary $g(\mathbf{x})$ is one-to-one with the likelihood ratio (LR) as $N \rightarrow \infty$

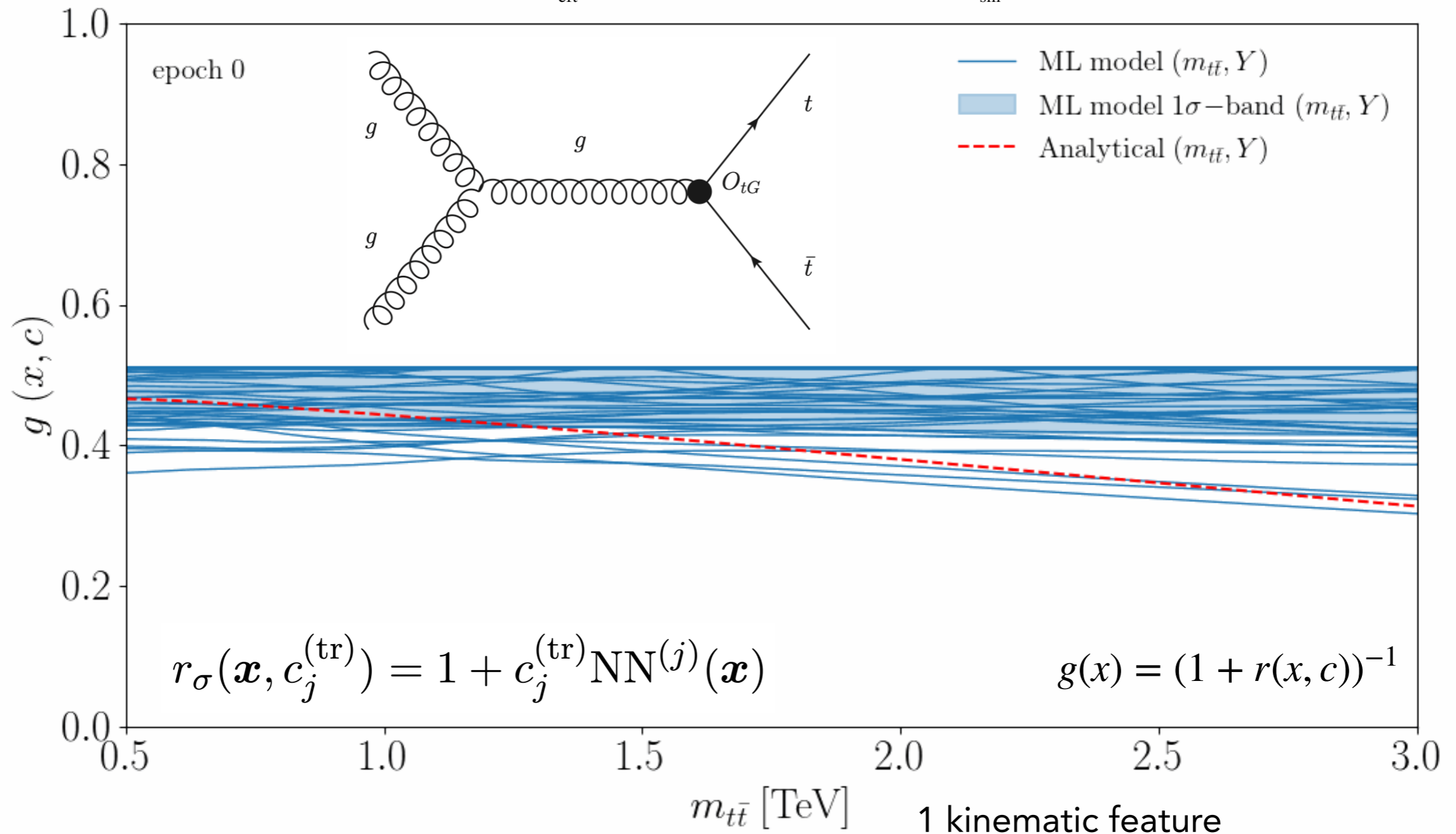
$$\frac{\delta L}{\delta g} = 0 \implies \hat{g}(\mathbf{x}) = \left(1 + \frac{f(\mathbf{x} | \text{EFT})}{f(\mathbf{x} | \text{SM})} \right)^{-1} \equiv \frac{1}{1 + r(\mathbf{x})}$$

Parameterise with NNs

- Parameterising $g(\mathbf{x})$ inside L with NNs lets us extract the likelihood (ratio) **implicitly**

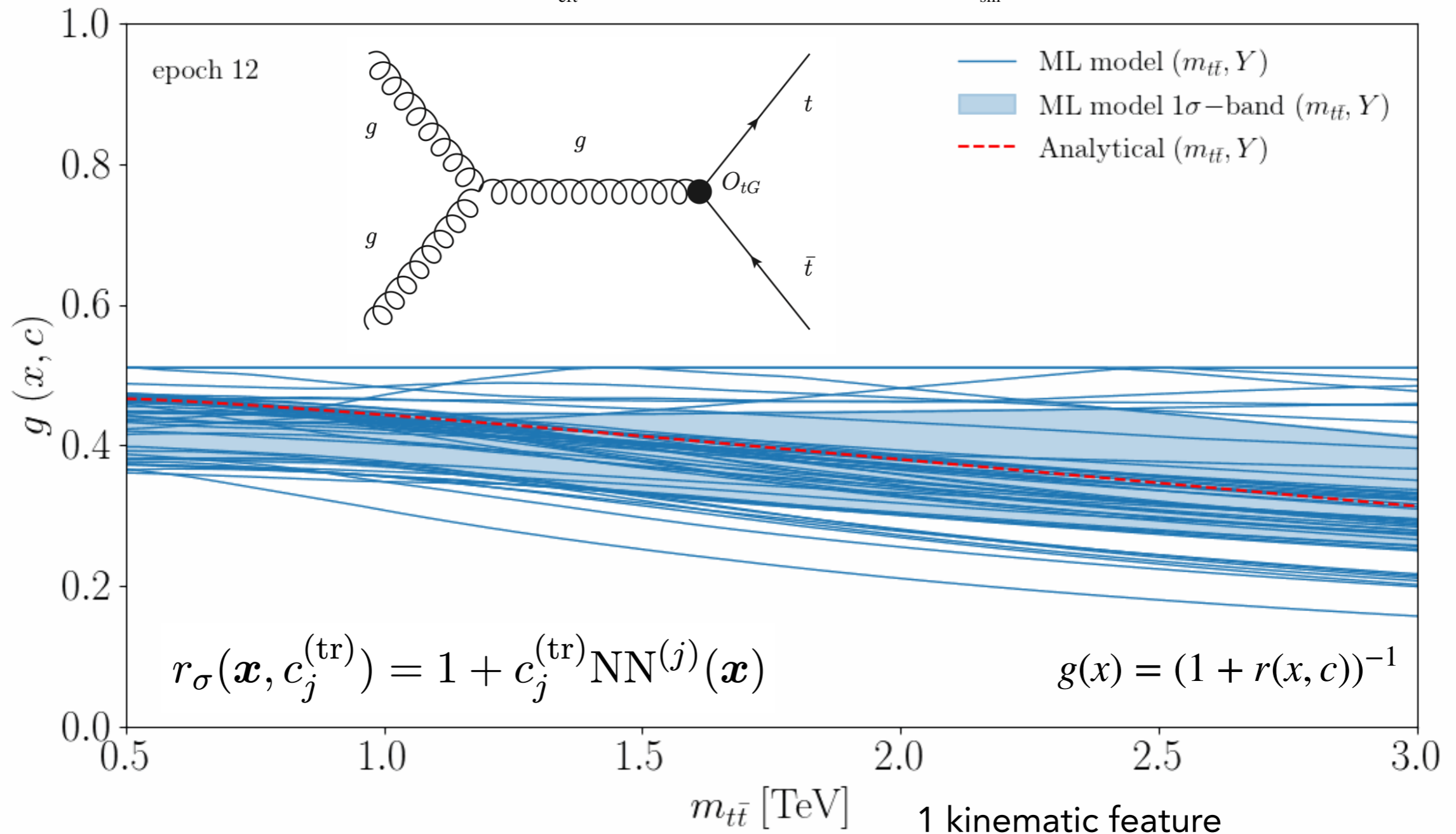
Likelihood learning in practice

$$L[g(\mathbf{x}, \mathbf{c})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{eff}}} w_e \log(1 - g(\mathbf{x}_e, \mathbf{c})) - \frac{1}{N} \sum_{e \in \mathcal{D}_{\text{sm}}} w_e \log g(\mathbf{x}_e, \mathbf{c})$$



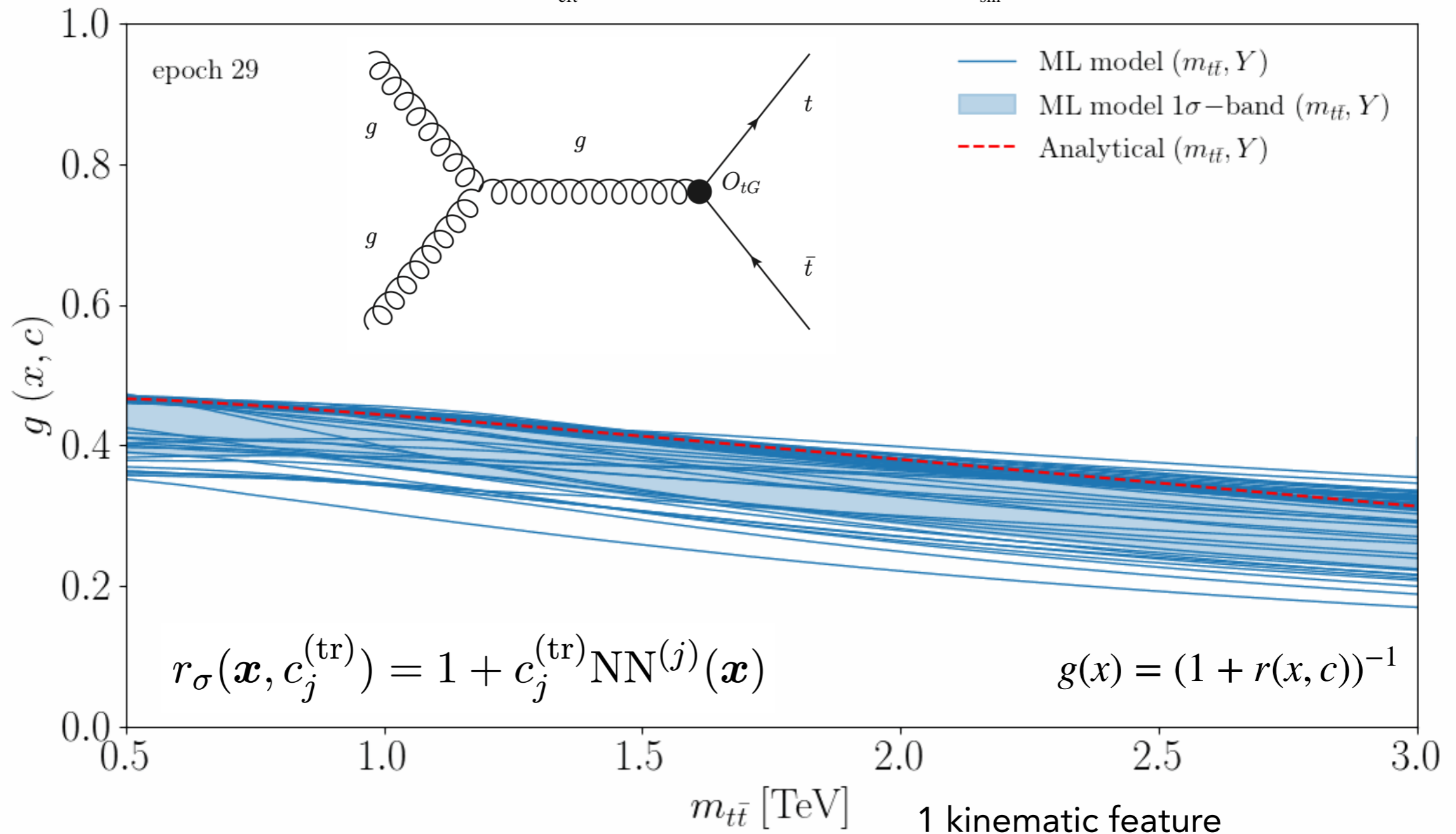
Likelihood learning in practice

$$L[g(\mathbf{x}, \mathbf{c})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{eff}}} w_e \log(1 - g(\mathbf{x}_e, \mathbf{c})) - \frac{1}{N} \sum_{e \in \mathcal{D}_{\text{sm}}} w_e \log g(\mathbf{x}_e, \mathbf{c})$$



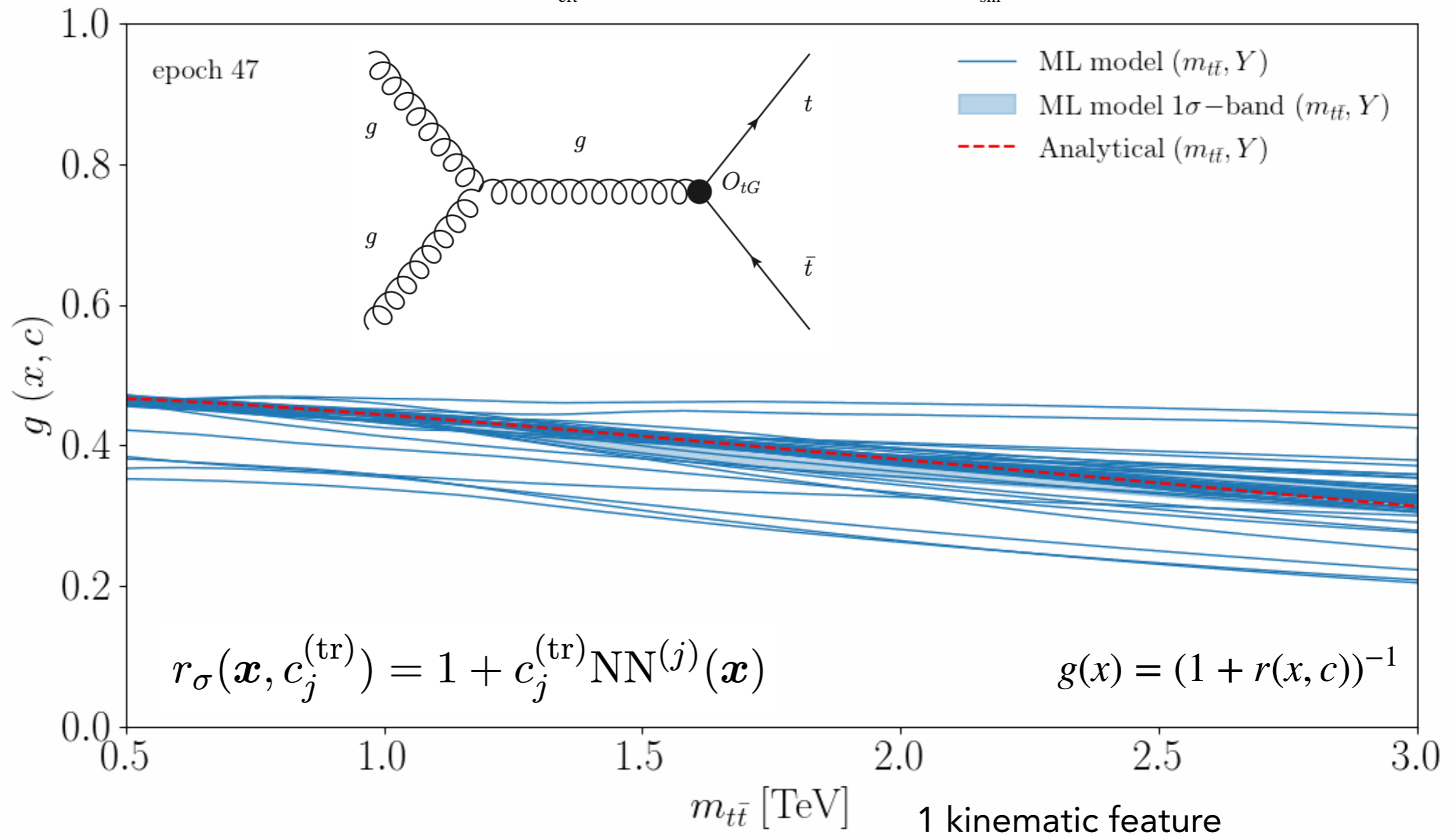
Likelihood learning in practice

$$L[g(\mathbf{x}, \mathbf{c})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{eff}}} w_e \log(1 - g(\mathbf{x}_e, \mathbf{c})) - \frac{1}{N} \sum_{e \in \mathcal{D}_{\text{sm}}} w_e \log g(\mathbf{x}_e, \mathbf{c})$$



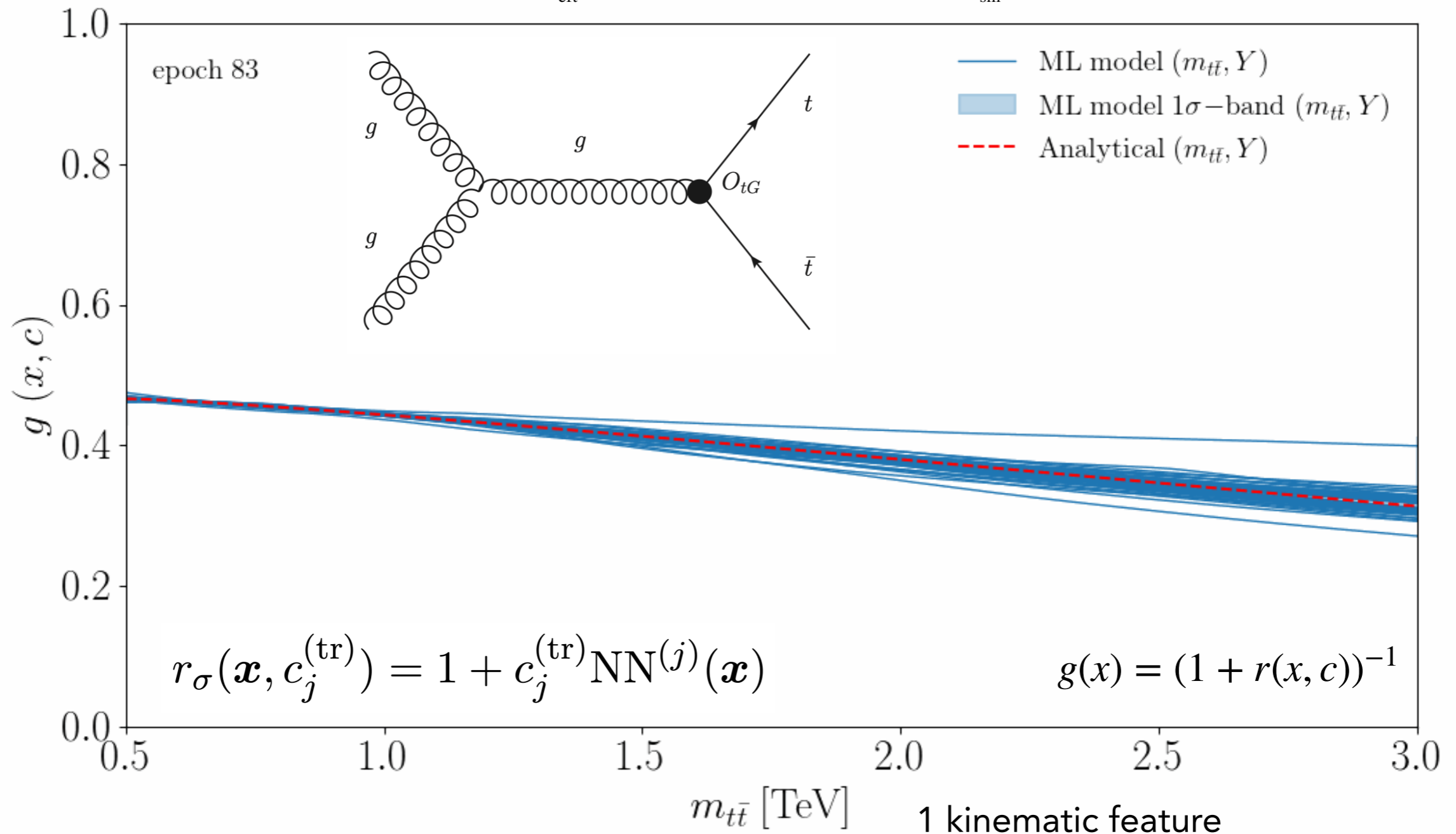
Likelihood learning in practice

$$L[g(\mathbf{x}, \mathbf{c})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{eff}}} w_e \log(1 - g(\mathbf{x}_e, \mathbf{c})) - \frac{1}{N} \sum_{e \in \mathcal{D}_{\text{sm}}} w_e \log g(\mathbf{x}_e, \mathbf{c})$$



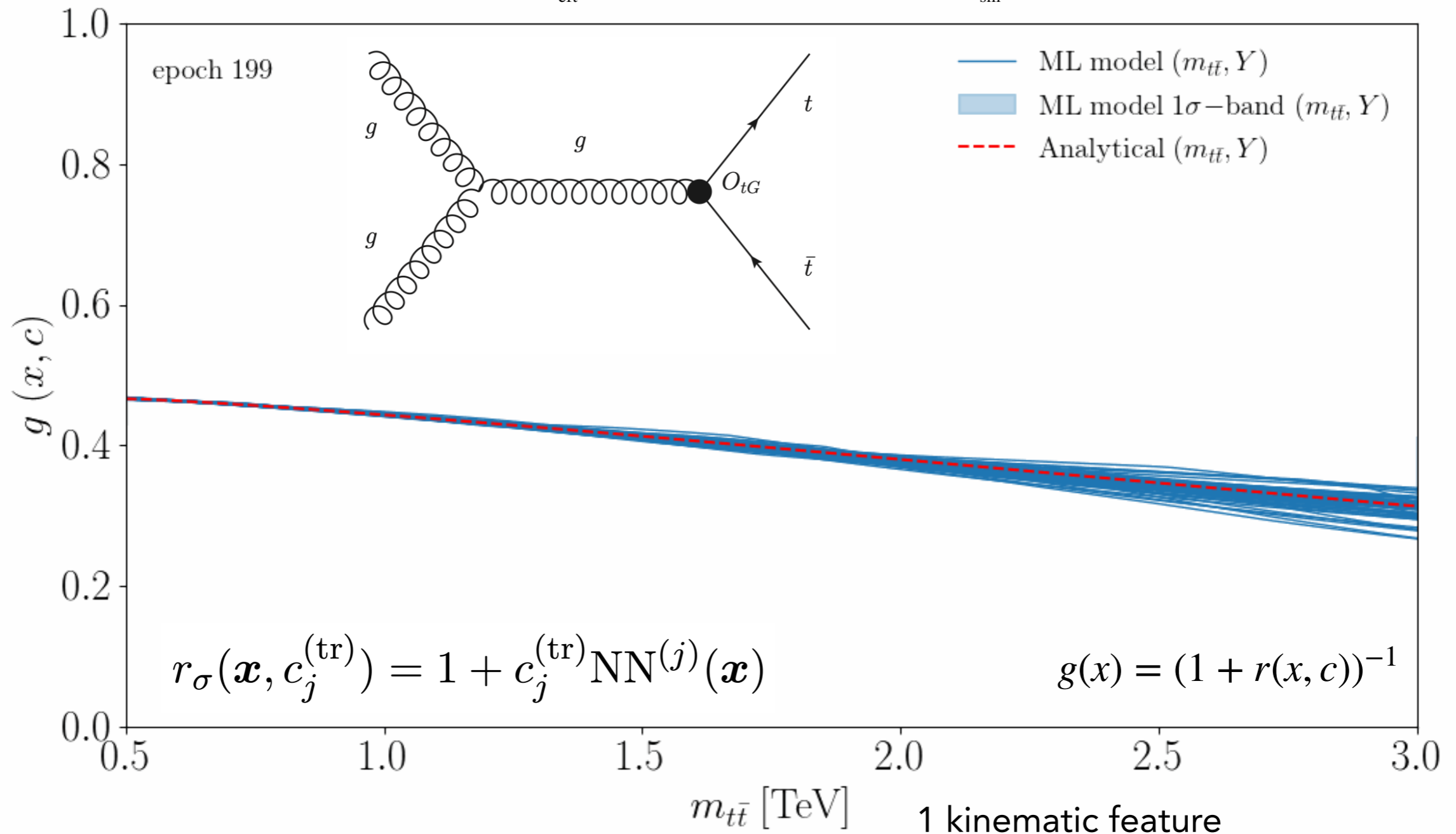
Likelihood learning in practice

$$L[g(\mathbf{x}, \mathbf{c})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{eff}}} w_e \log(1 - g(\mathbf{x}_e, \mathbf{c})) - \frac{1}{N} \sum_{e \in \mathcal{D}_{\text{sm}}} w_e \log g(\mathbf{x}_e, \mathbf{c})$$



Likelihood learning in practice

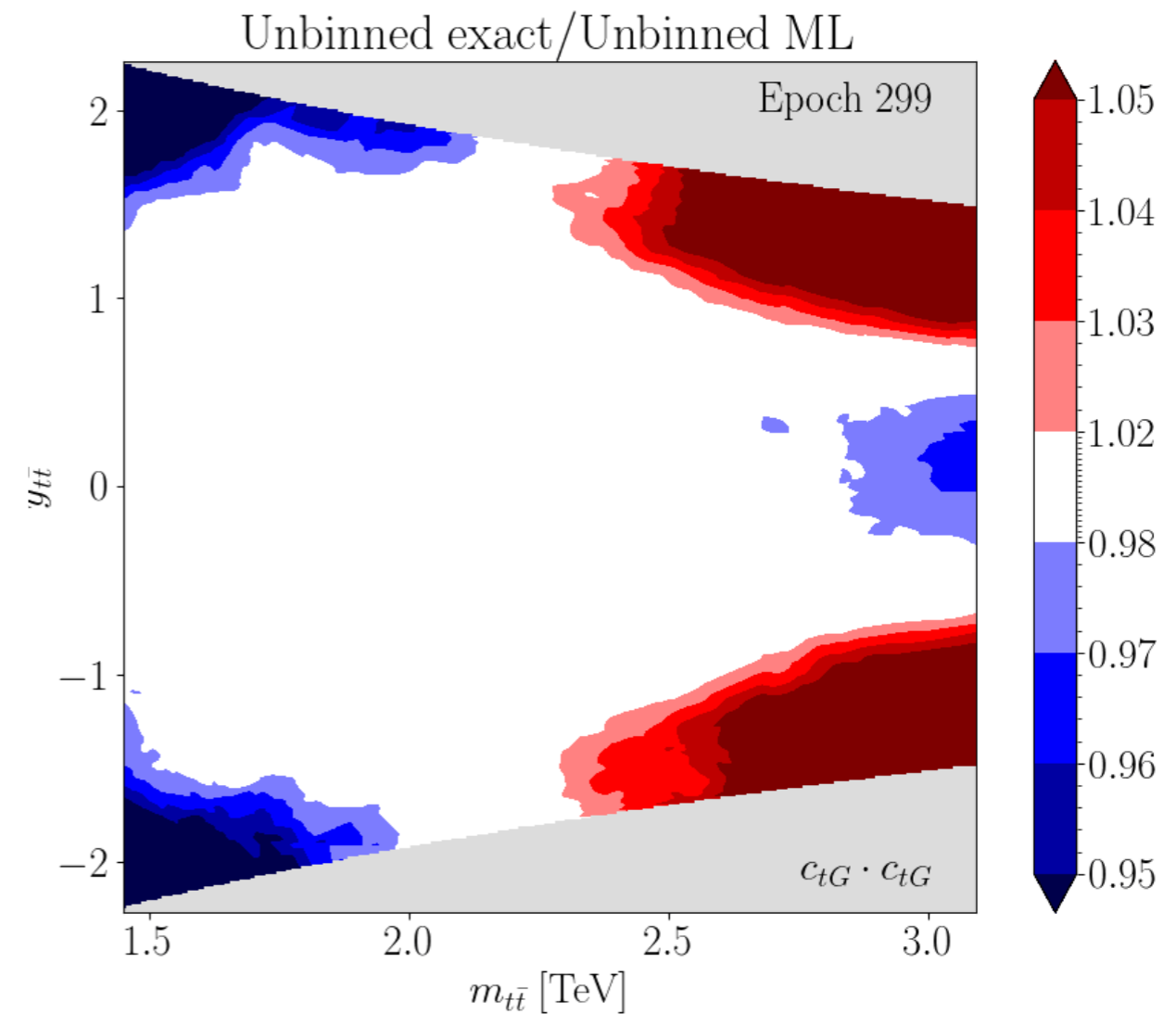
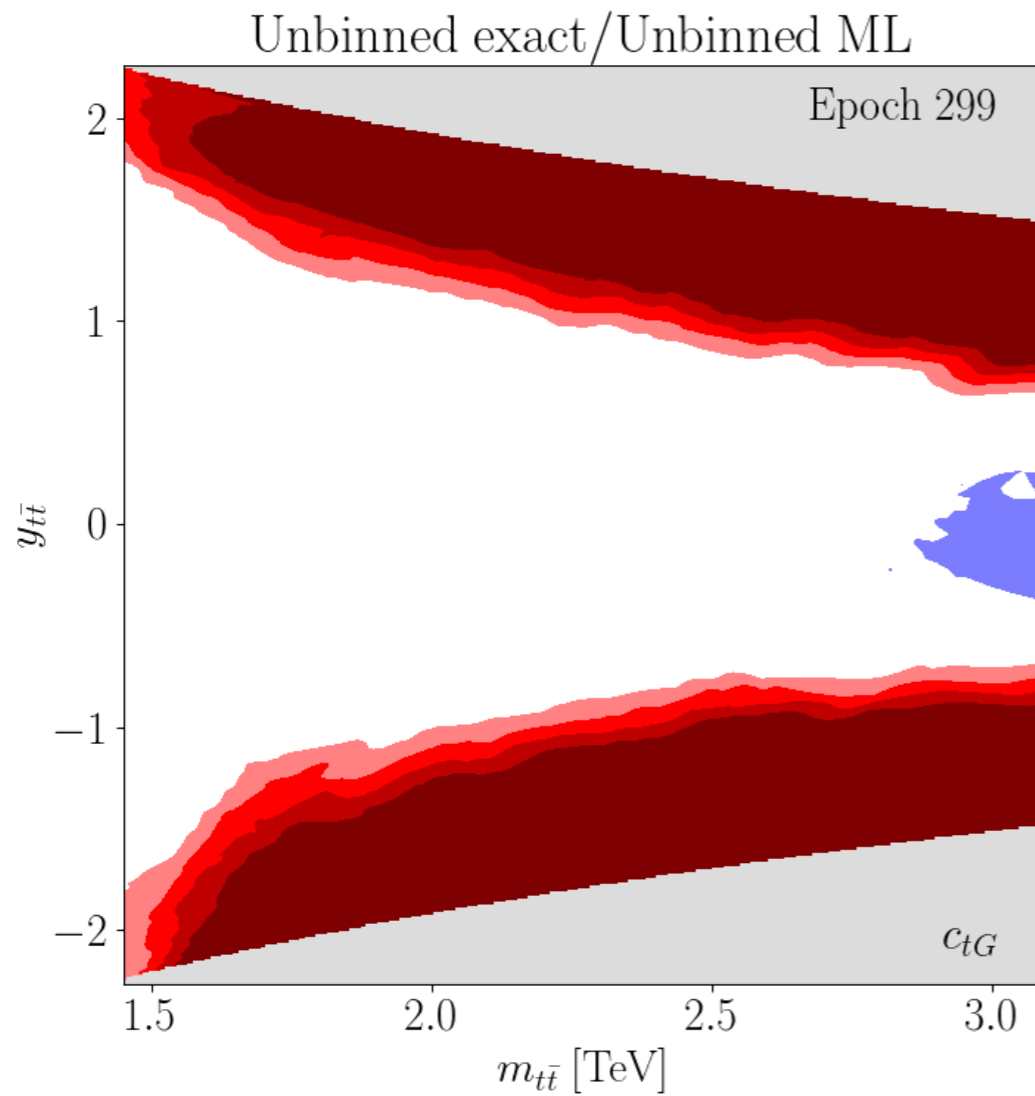
$$L[g(\mathbf{x}, \mathbf{c})] = -\frac{1}{N} \sum_{e \in \mathcal{D}_{\text{eff}}} w_e \log(1 - g(\mathbf{x}_e, \mathbf{c})) - \frac{1}{N} \sum_{e \in \mathcal{D}_{\text{sm}}} w_e \log g(\mathbf{x}_e, \mathbf{c})$$



Likelihood learning in practice

2 kinematic features

$$r_{\sigma}(\mathbf{x}, c_j^{(\text{tr})}) = 1 + c_j^{(\text{tr})} \text{NN}^{(j)}(\mathbf{x})$$



The ML4EFT framework

`pip install ml4eft`

<https://lhcfitnikhef.github.io/ML4EFT>

2211.02058 R. Gomez Ambrosio, JtH, M. Madigan, J. Rojo, V.Sanz

Open-source NN-based python framework for the integration of unbinned multivariate observables into global SMEFT fits

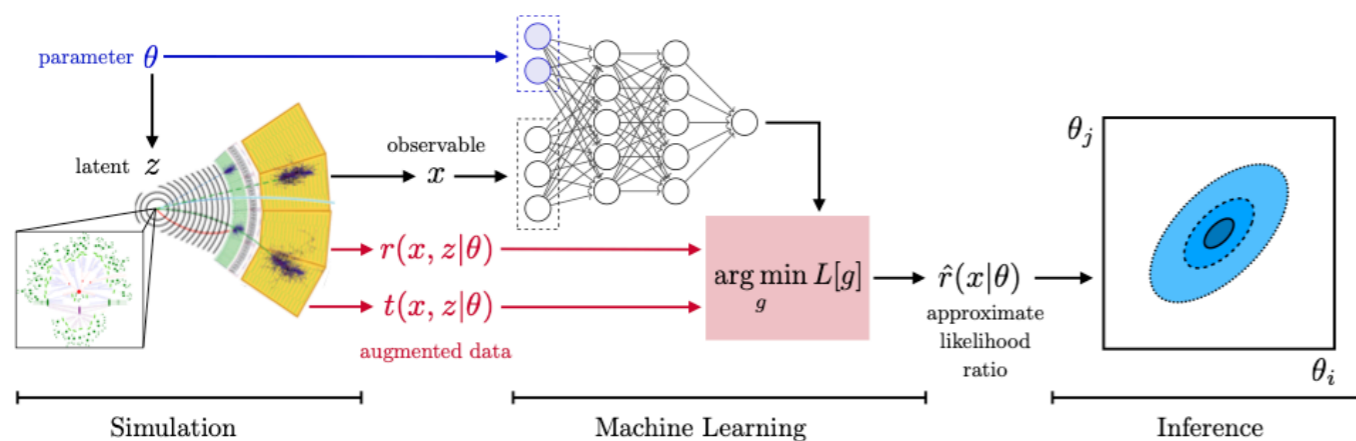
- ▶ **Goal:** provide optimal constraints on the SMEFT
- ▶ **Diagnostic tool:** what is the information loss incurred by a particular choice of bins?
- ▶ **Projections:** how will SMEFT constraints improve if unbinned data are made available?

The screenshot displays the documentation for the `ml4eft.core.classifier.Fitter` class. On the left, a navigation menu lists various parts of the documentation, including 'Installation', 'Tutorial', and a tree view of the package structure under 'ml4eft'. The main content area shows the class definition: `class ml4eft.core.classifier.Fitter(json_path, mc_run, c_name, output_dir, print_log=False)`. It lists the base class as `object` and identifies it as a 'Training class'. The `__init__` method is highlighted, with its signature and a link to the source code. Below this, the 'Fitter constructor' section lists parameters: `json_path` (Path to json run card), `mc_run` (Replica number), `c_name` (EFT coefficient for which to learn the ratio function), `output_dir` (Path to where the models should be stored), and `print_log` (Set to true to print training progress to stdout, otherwise it prints to a log file only). A 'Methods' section lists several methods: `__init__` (Fitter constructor), `load_data()` (Constructs training and validation sets), `loss_fn(outputs, labels, w_e)` (Loss function), `train_classifier(data_train, data_val)` (Starts the training of the binary classifier), `training_loop(optimizer, train_loader, ...)` (Optimize the classifier with `optimizer` on the training data set `train_loader`), and `weight_reset(m)` (Reset the weights and biases associated with the model `m`).

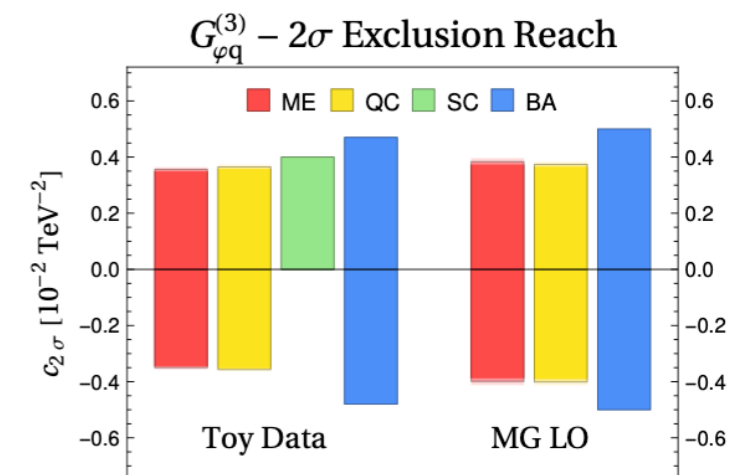
Modular structure, easy to maintain, well documented

Related work

- ▶ Madminer series (J.Brehmer, K.Cranmer, G.Louppe et al.) [1907.10621, 1805.00020, ...]
- ▶ Parameterized classifiers for SMEFT (A. Glioti et al.) [2007.10356]
- ▶ Learning the EFT likelihood with tree boosting (R. Schöfbeck et al) [2205.12976]
- ▶ Back to the Formula (A. Butter, T. Plehn et al) [2109.10414]
- ▶ Boosted likelihood learning with event reweighing (A. Glioti et al) [2308.05704]
- ▶ Designing Observables for Measurements with Deep Learning (O.Long, B. Nachman) [2310.08717]



[2010.06439]



[2007.10356]

Anticipating global fits

- Global EFT fits typically feature ~ 50 WCs and thus efficient scaling with the number of WCs becomes essential
- Solution: learn the coefficient functions separately and combine afterwards

$$r(\mathbf{x}, \mathbf{c}) \equiv \frac{d\sigma(\mathbf{x}, \mathbf{c})}{d\sigma(\mathbf{x}, \mathbf{0})} = 1 + \sum_{j=1}^{n_{\text{eft}}} r^{(j)}(\mathbf{x}) c_j + \sum_{j=1}^{n_{\text{eft}}} \sum_{k \geq j}^{n_{\text{eft}}} r^{(j,k)}(\mathbf{x}) c_j c_k$$

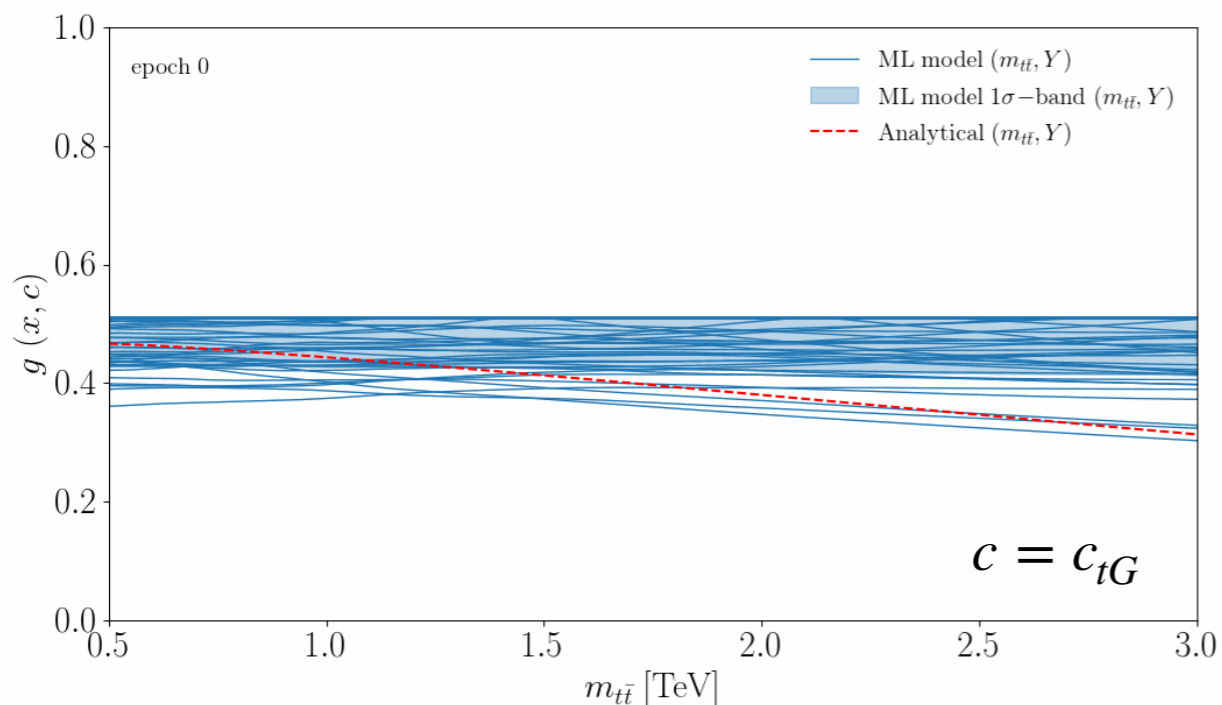
Example: to learn a single $r^{(j)}$, generate \mathcal{D}_{sm} and \mathcal{D}_{eft} at c_j up to $\mathcal{O}(\Lambda^{-2})$.
Then $r(\mathbf{x}, \mathbf{c}) = 1 + r^{(j)}(\mathbf{x}) c_j^{(\text{tr})}$ and training means

$$g(\mathbf{x}, c_j^{(\text{tr})}) = \left(1 + \left[1 + c_j^{(\text{tr})} \cdot \text{NN}^{(j)}(\mathbf{x}) \right] \right)^{-1} \quad \text{NN}^{(j)}(\mathbf{x}) \rightarrow r^{(j)}(\mathbf{x})$$

Uncertainty treatment

- ▶ Stick to a regime in which statistical uncertainties dominate over systematics
- ▶ Finite training data makes one subject to methodological uncertainties
- ▶ Solution: propagate uncertainties to the space of models by training multiple replicas

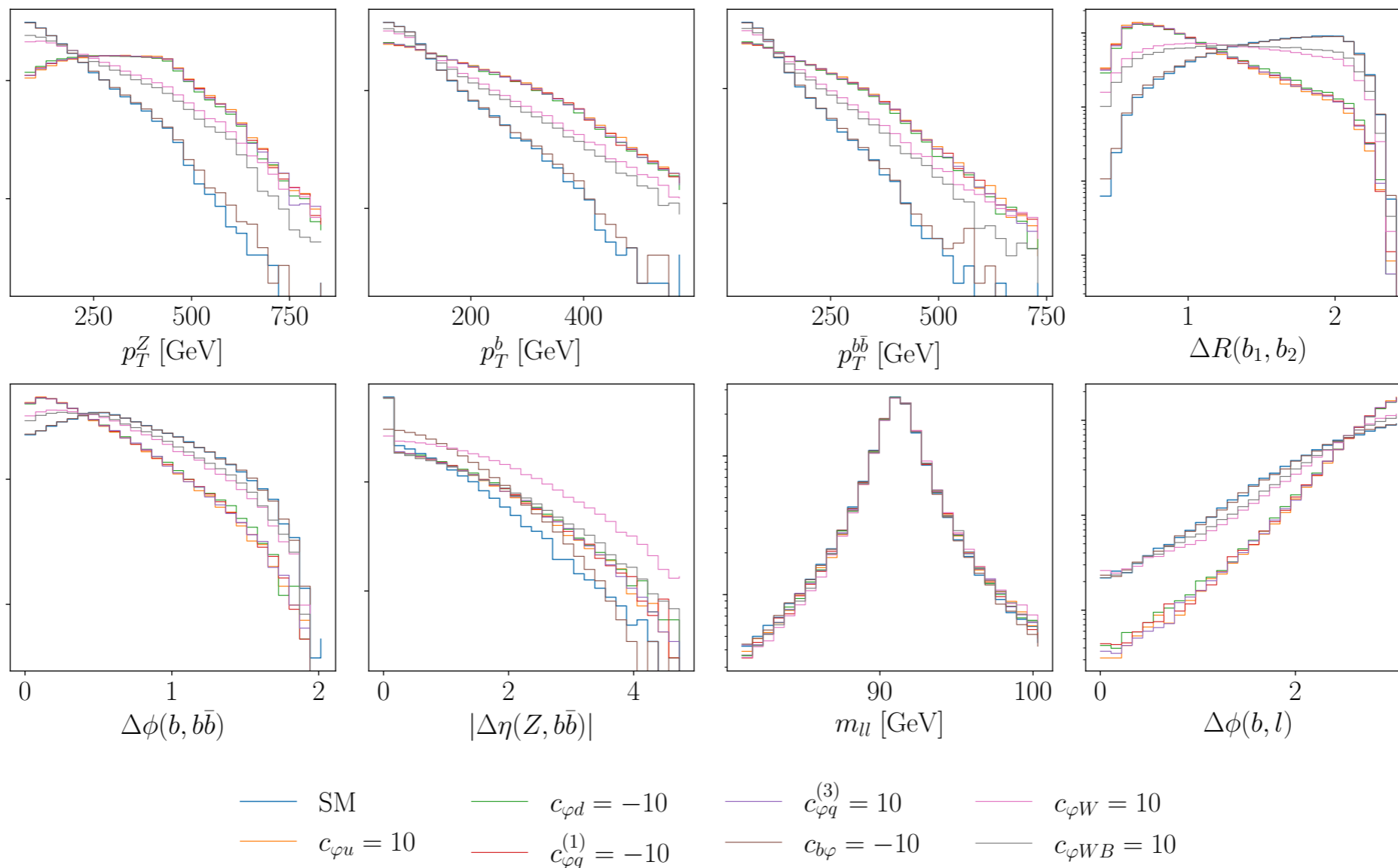
$$\hat{r}^{(i)}(\mathbf{x}, \mathbf{c}) \equiv 1 + \sum_{j=1}^{n_{\text{eff}}} \text{NN}_i^{(j)}(\mathbf{x})c_j + \sum_{j=1}^{n_{\text{eff}}} \sum_{k \geq j}^{n_{\text{eff}}} \text{NN}_i^{(j,k)}(\mathbf{x})c_j c_k, \quad i = 1, \dots, N_{\text{rep}}$$



Process	N_{rep}	\tilde{N}_{ev} (per replica)	N_{nn}	#trainings
$pp \rightarrow t\bar{t}$	50	10^5	4	200
$pp \rightarrow t\bar{t} \rightarrow b\bar{b}l^+\ell^-\nu_\ell\bar{\nu}_\ell$	25	10^5	40	1000

Let's go multivariate

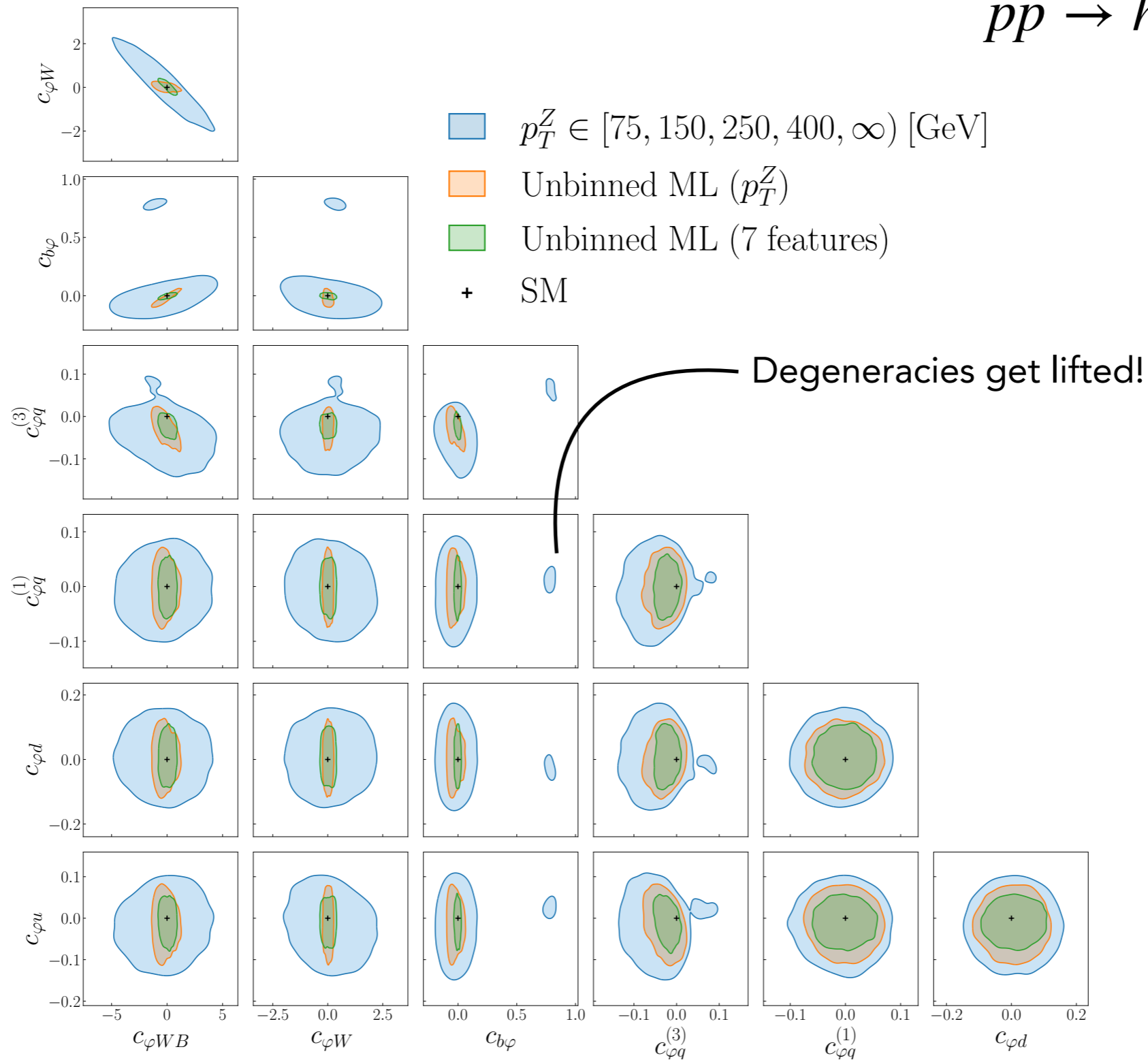
- $pp \rightarrow t\bar{t} \rightarrow b\bar{b}\ell^+\ell^-\nu_\ell\bar{\nu}_\ell$: 18 features, 8 EFT coefficients
- $pp \rightarrow hZ \rightarrow b\bar{b}\ell^+\ell^-$: 7 features, 7 EFT coefficients



Results: Higgs + Z associated production

Marginalised 95 % C.L. intervals, $\mathcal{O}(\Lambda^{-4})$ at $\mathcal{L} = 300 \text{ fb}^{-1}$

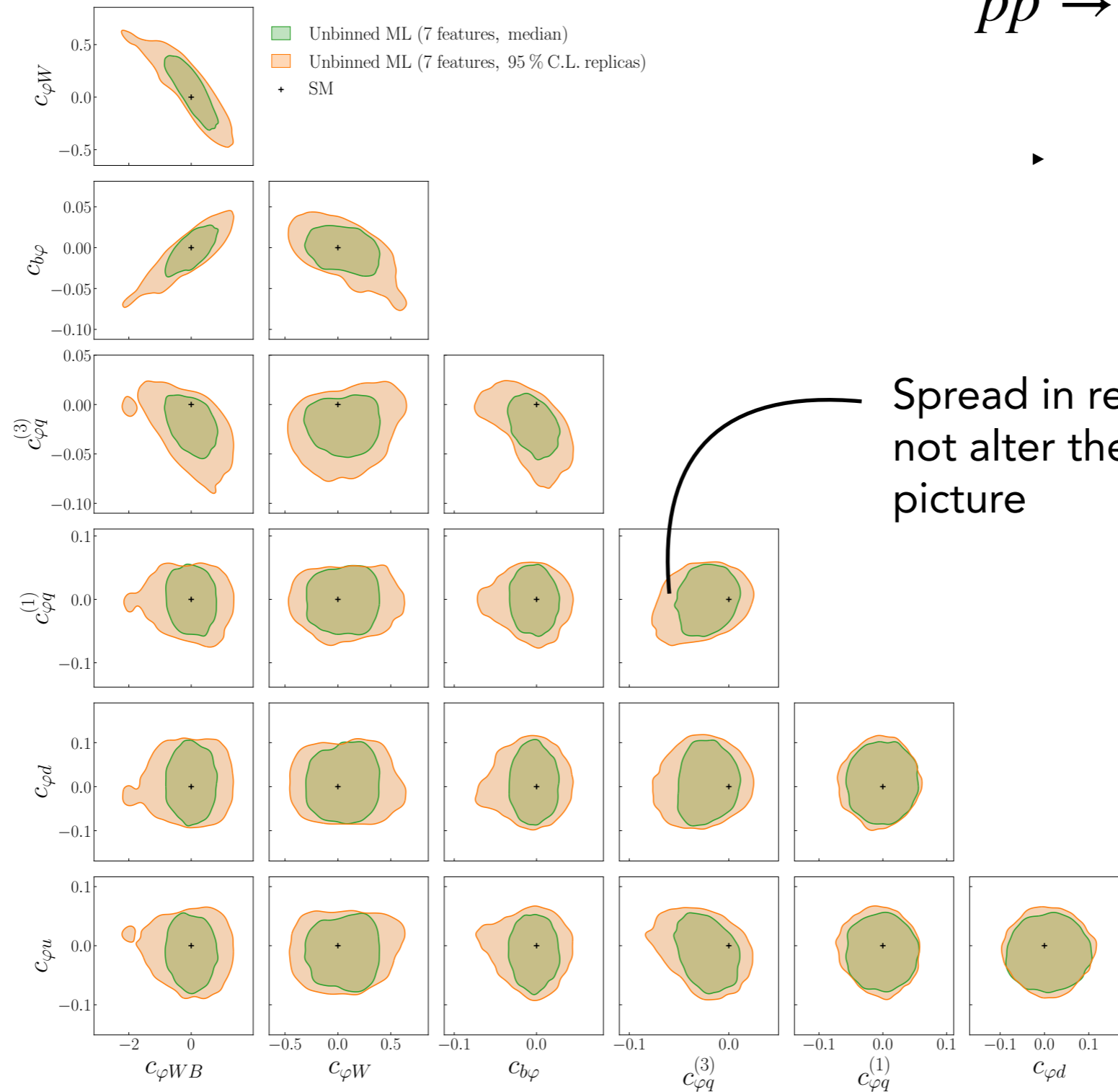
$$pp \rightarrow hZ \rightarrow b\bar{b}\ell^+\ell^-$$



Results: Higgs + Z associated production

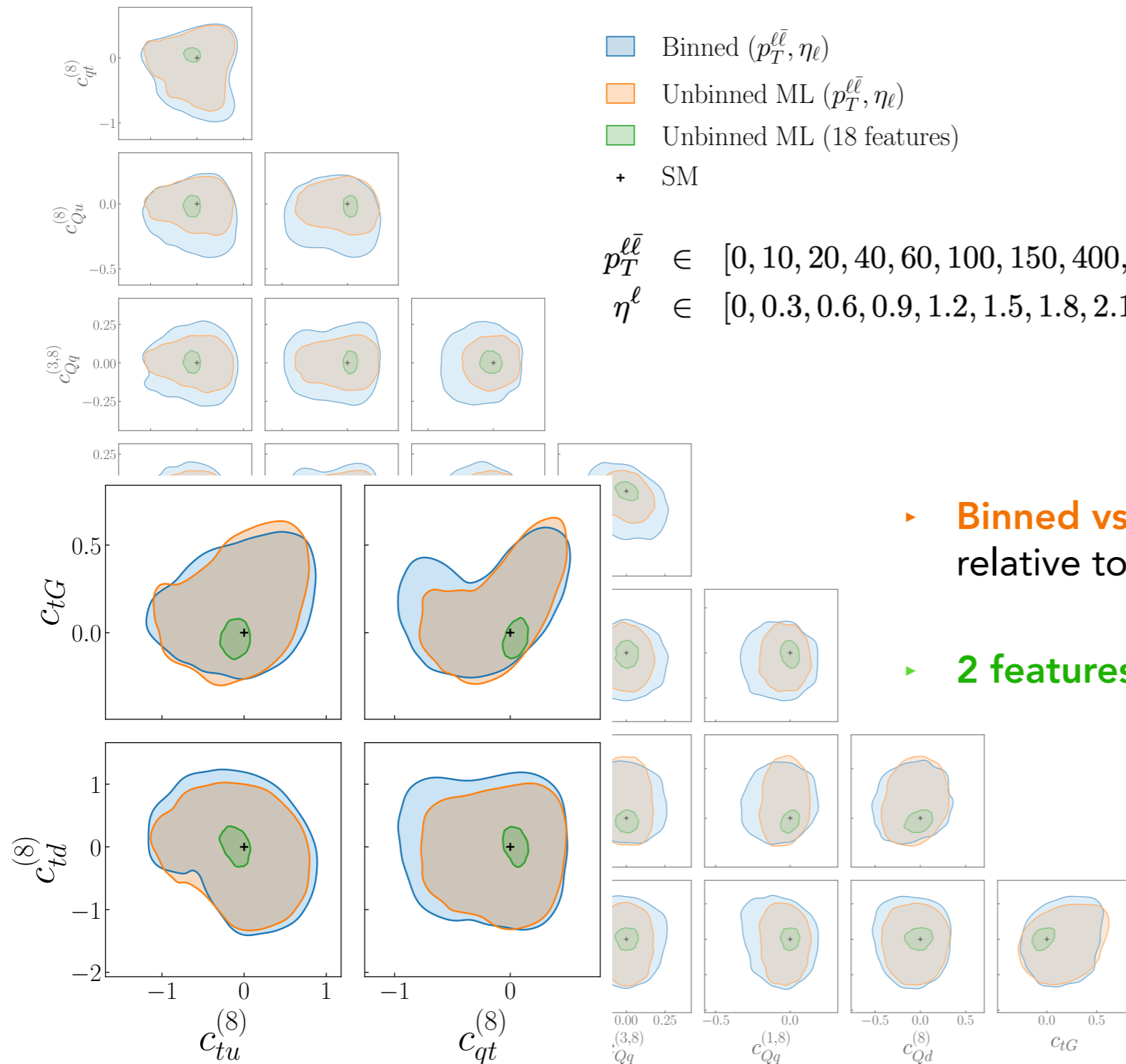
Marginalised 95 % C.L. intervals, $\mathcal{O}(\Lambda^{-4})$ at $\mathcal{L} = 300 \text{ fb}^{-1}$

$$pp \rightarrow hZ \rightarrow b\bar{b}\ell^+\ell^-$$



Results: unbinned observables in the top sector

Marginalised 95 % C.L. intervals, $\mathcal{O}(\Lambda^{-4})$ at $\mathcal{L} = 300 \text{ fb}^{-1}$

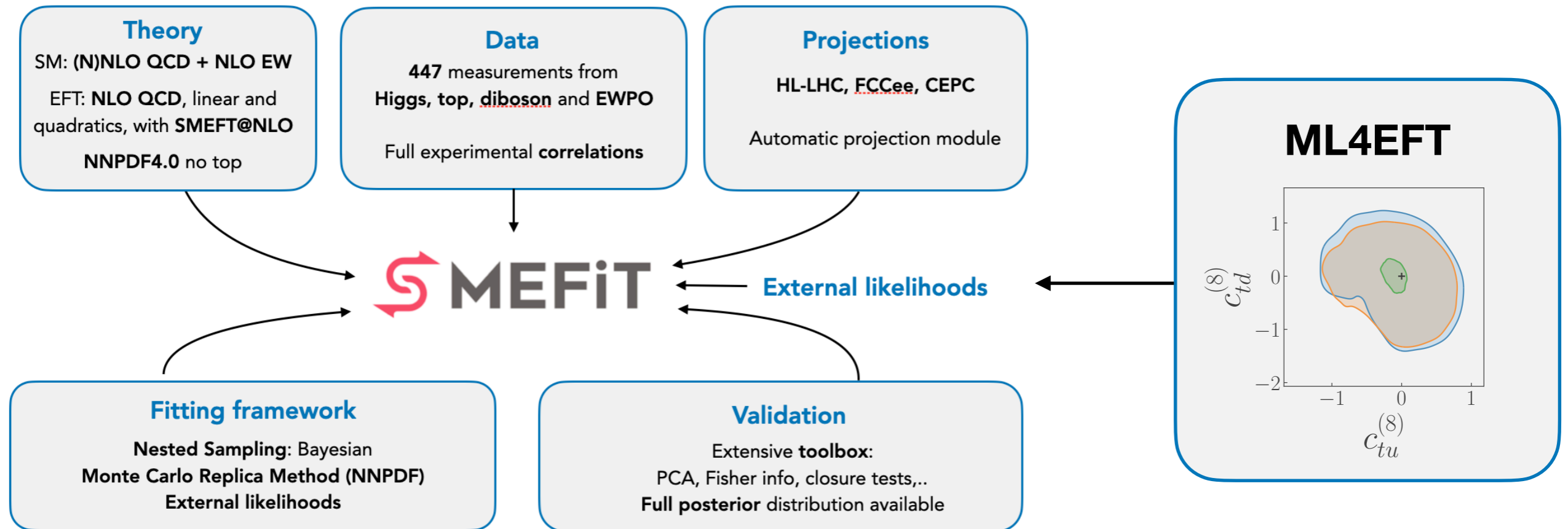


- ▶ **Binned vs unbinned** in $(p_T^{\ell\bar{\ell}}, \eta_\ell)$ small improvement relative to binned setup
- ▶ **2 features vs 18 features:** big increase in sensitivity

ML4EFT + SMEFIT

A combined framework

Best of two worlds?



- The ultimate global EFT fit combines **binned** and **multivariate unbinned ML** observables
- We need a framework that connects them

A combined framework

The SMEFiT (binned) likelihood



$$-2 \log \mathcal{L}(c) = \frac{1}{n_{\text{dat}}} \sum_{i,j=1}^{n_{\text{dat}}} \left(\sigma_{i,\text{SMEFT}}(c) - \sigma_{i,\text{exp}} \right) (\text{cov}^{-1})_{ij} \left(\sigma_{j,\text{SMEFT}}(c) - \sigma_{j,\text{exp}} \right)$$

And the multivariate unbinned ML likelihood (ratio)



$$\mathcal{L}(c) = \frac{\nu_{\text{tot}}(c)^{N_{\text{ev}}}}{N_{\text{ev}}!} e^{-\nu_{\text{tot}}(c)} \prod_{i=1}^{N_{\text{ev}}} f_{\sigma}(\mathbf{x}_i, c)$$

Together form a powerful overall likelihood



$$\log \mathcal{L}(c) = \sum_{k=1}^{N_D^{(\text{unbinned})}} \log \mathcal{L}_k^{\text{unbinned}}(c) + \sum_{k=1}^{N_D^{(\text{binned})}} \log \mathcal{L}_k^{\text{binned}}(c)$$

A toy setup



- ATLAS_CMS_SSinc_RunI
- ATLAS_SSinc_RunII
- CMS_SSinc_RunII

- ATLAS_WW_13TeV_2016_memu
- ATLAS_WZ_13TeV_2016_mTWZ
- CMS_WZ_13TeV_2016_pTZ
- CMS_WZ_13TeV_2022_pTZ

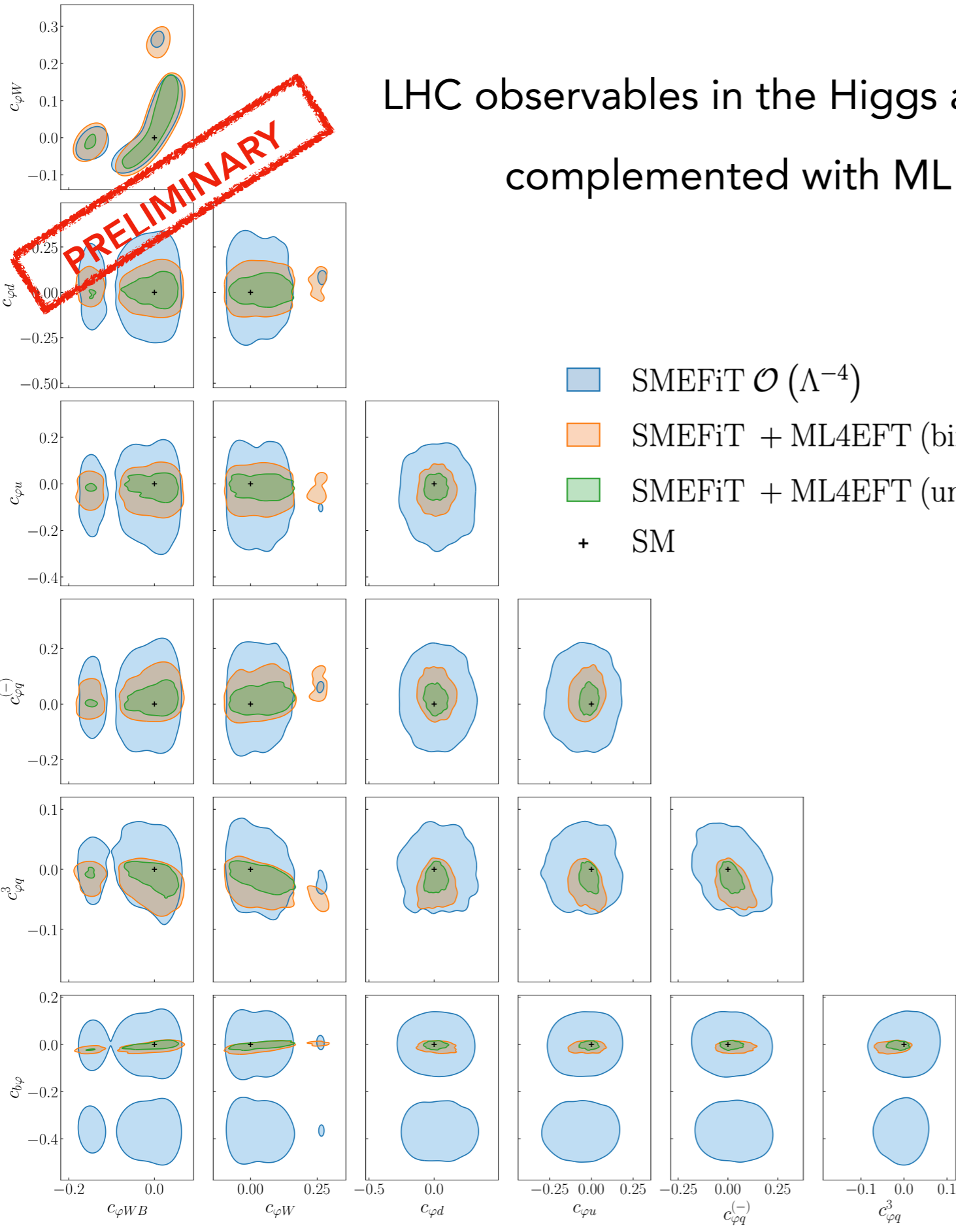
- ATLAS_STXS_runII_13TeV
- ATLAS_WH_Hbb_13TeV
- ATLAS_ZH_Hbb_13TeV
- ATLAS_ggF_13TeV_2015
- CMS_ggF_aa_13TeV
- ATLAS_ggF_ZZ_13TeV
- CMS_H_13TeV_2015_pTH

+

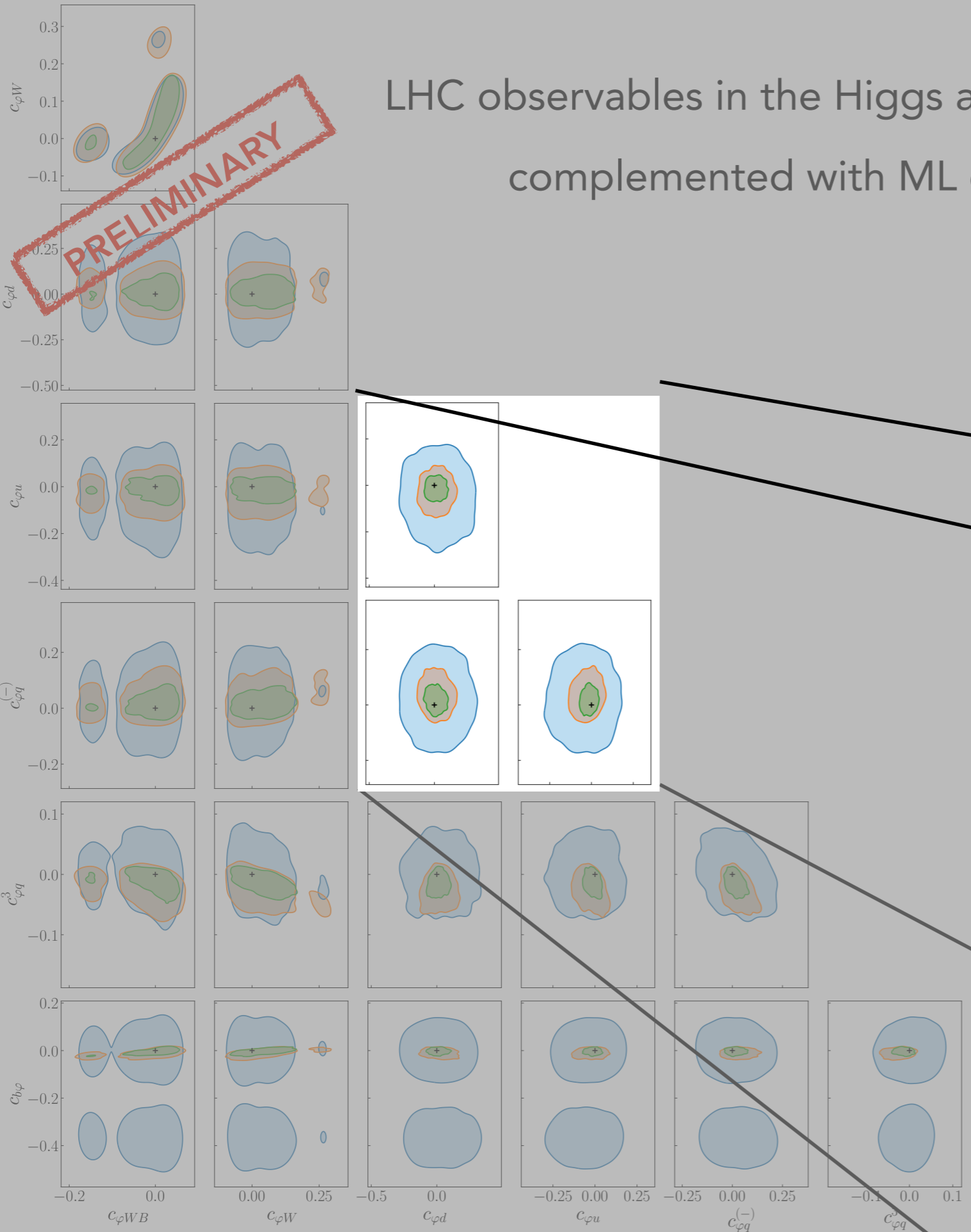
ML4EFT Run III projections

ATLAS_ZH_Hbb_13TeV

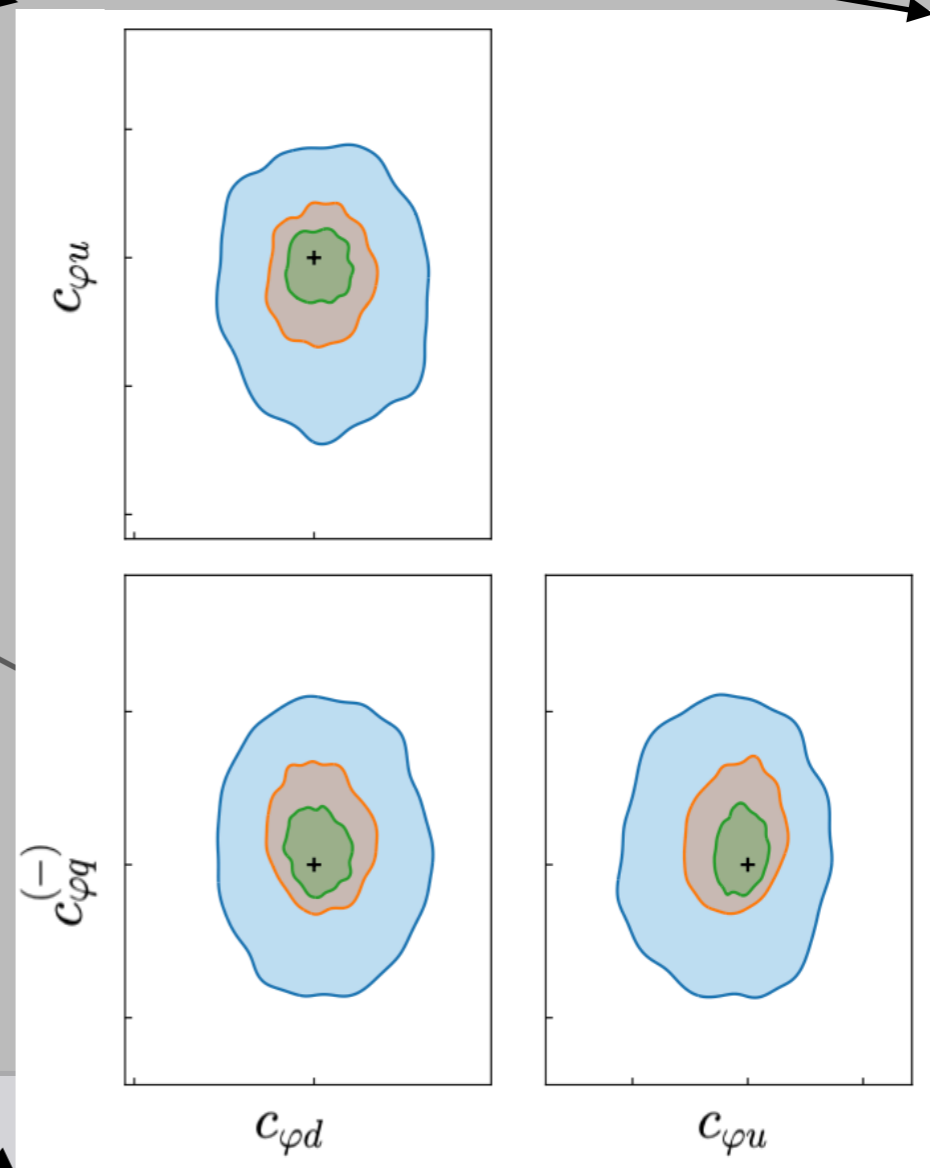
LHC observables in the Higgs and diboson sector... complemented with ML observables



LHC observables in the Higgs and diboson sector.... complemented with ML observables

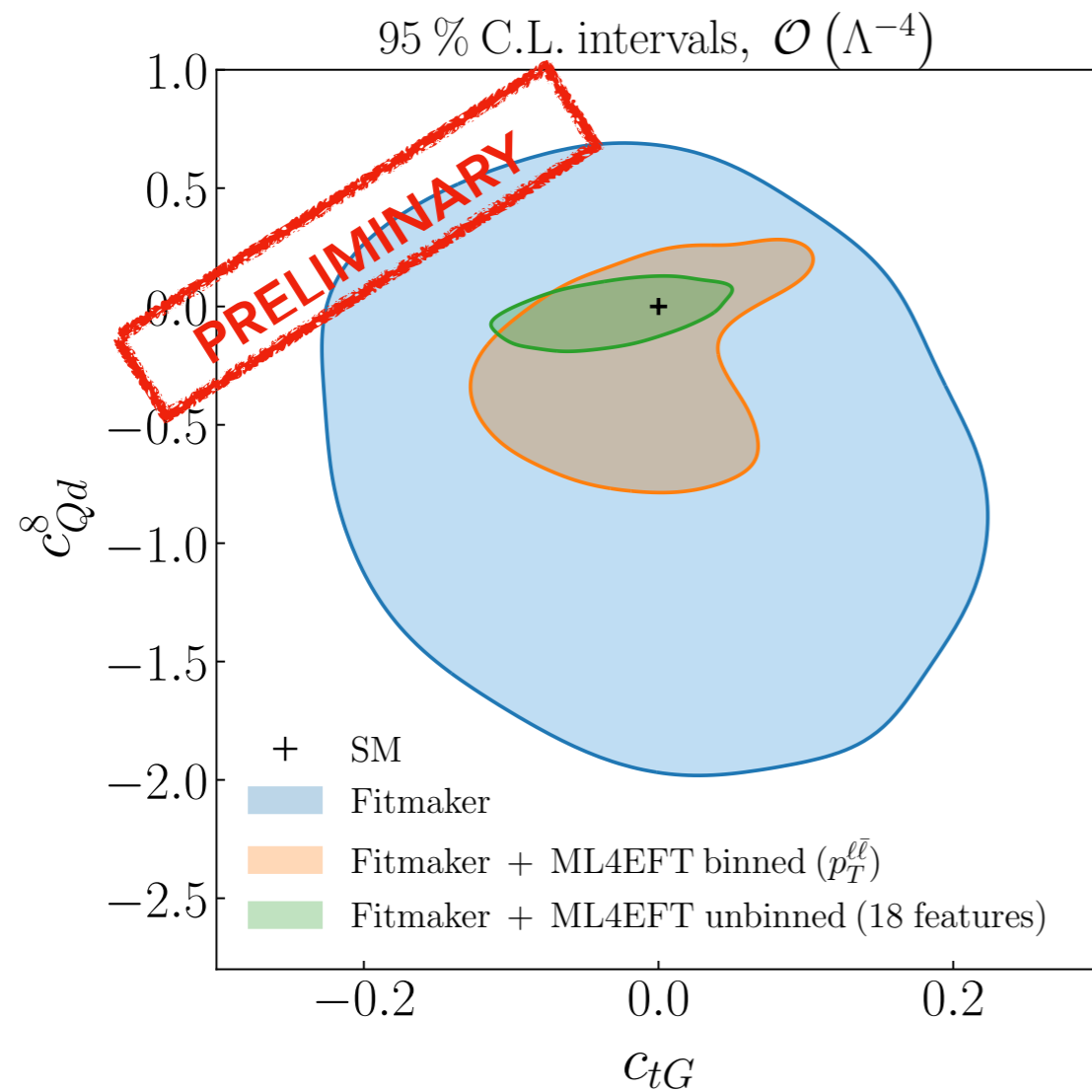


- SMEFiT $\mathcal{O}(\Lambda^{-4})$
- SMEFiT + ML4EFT (binned, p_T^Z) $\mathcal{O}(\Lambda^{-4})$
- SMEFiT + ML4EFT (unbinned, 7 features) $\mathcal{O}(\Lambda^{-4})$
- +
- SM



A combined framework

- Similar picture in the top sector



Summary

- State of the art global SMEFT fit with the latest run II results
- ML4EFT integrates unbinned multivariate observables into global SMEFT fits with a faithful uncertainty estimate through the replica method
- Case study in the Higgs and top sector
- Global EFT fits show increased sensitivity to unbinned observables
- Please visit **ML4EFT** on GitHub (documentation + **tutorial**)

lhcfiteknikhef.github.io/ML4EFT

Thank you!